

# GENERATIVE REPRESENTATIONAL INSTRUCTION TUNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

All text-based language problems can be reduced to either generation or embedding. Current models only perform well at one or the other. We introduce generative representational instruction tuning (GRIT) whereby a large language model is trained to handle both generative and embedding tasks by distinguishing between them through instructions. Compared to other open models, our resulting GRITLM 7B is among the top models on the Massive Text Embedding Benchmark (MTEB) and outperforms various models up to its size on a range of generative tasks. By scaling up further, GRITLM 8x7B achieves even stronger generative performance while still being among the best embedding models. Notably, we find that GRIT matches training on only generative or embedding data, thus we can unify both at no performance loss. Among other benefits, the unification via GRIT speeds up Retrieval-Augmented Generation (RAG) by >60% for long documents, by no longer requiring separate retrieval and generation models. Models, code, etc. will be made freely available.

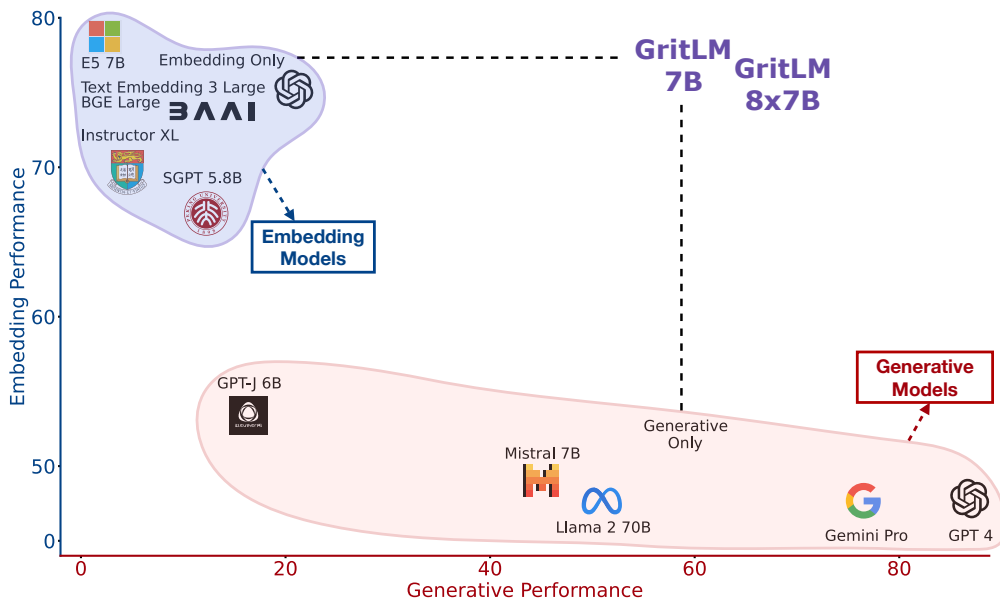


Figure 1: **Performance of various models on text representation (embedding) and generation tasks.** GRITLM is the first model to perform strongly at both types of tasks simultaneously.

## 1 INTRODUCTION

Creating a single general model that performs well at a wide range of tasks has been a long-standing goal of the field of artificial intelligence (Kaiser et al., 2017; Jaegle et al., 2021; Cho et al., 2021; Reed et al., 2022; Singh et al., 2022). Recently, large language models (LLMs) have emerged as a

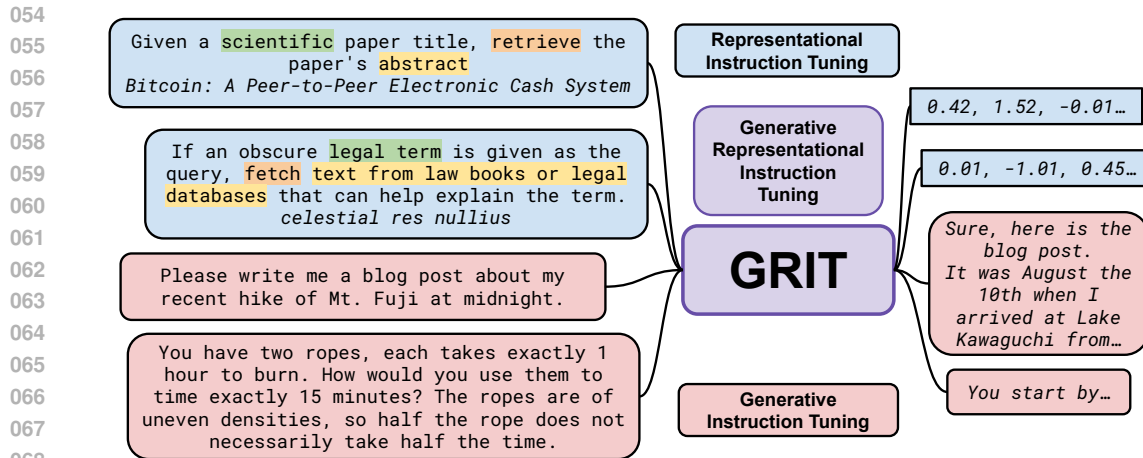


Figure 2: **GRIT**. The same model handles both text representation and generation tasks based on a given instruction. For representation tasks, instructions ideally contain target **domain**, **intent**, and **unit** (Asai et al., 2022). The representation is a numeric tensor, while the generative output is text.

promising direction for a single multi-task model (Radford et al., 2019; Brown et al., 2020). Prior work has argued that all text-based language problems can be reduced to generation and thus handled by a single LLM (Raffel et al., 2023; Du et al., 2021).

However, tasks that use embeddings, such as clustering or retrieval (Muennighoff et al., 2023c), have largely been ignored from this perspective. Today, text embeddings power many real-world applications ranging from search engines to user-facing chatbots (Huang et al., 2020; Su et al., 2017). While integrating text embeddings into the generative paradigm is possible by generating a sequence of numbers to form the embedding tensor, it becomes impractical due to the high dimensionality and precision requirements of embeddings. Thus, it is more common and much easier to use the hidden state of the model as the embedding representation, which is already a numeric tensor (Muennighoff, 2022; Wang & Kuo, 2020; Morris et al., 2023). However, for current generative models this leads to poor performance. For example, while the T5 model (Raffel et al., 2023; Sanh et al., 2022) can handle any generative task in a sequence-to-sequence fashion, it requires finetuning to make its hidden state useful for text embedding (Ni et al., 2021a;b) during which it loses its generative capabilities.

We introduce GRIT (generative representational instruction tuning) which unifies embedding and generative tasks, leading to a model that excels at both tasks as shown in Figure 1. Figure 2 depicts how GRIT combines two previously disjoint training paradigms: (1) *Generative instruction tuning*, whereby the model is trained to respond to instructions by generating an answer (Wei et al., 2022; Sanh et al., 2022); and (2) *Representational instruction tuning*, whereby the model is trained to represent a provided input according to an instruction (Su et al., 2023; Asai et al., 2022). Via the instructions and separate loss functions the model learns to differentiate the two streams. We test our approach on models with up to 47B parameters. This unification via GRIT leads to three advantages:

**a) Performance:** Our unified model matches the performance of embedding-only and generative-only variants, even outperforming them on some tasks. At 7B parameters, GRITLM is among the best models on the Massive Text Embedding Benchmark (Muennighoff et al., 2023c) and at the same time outperforms some larger models on generative tasks, such as Llama 2 70B. By scaling further, GRITLM 8x7B achieves even stronger generative performance, while only using 13B parameters at inference due to its MoE architecture (Jiang et al., 2024). Further, as our models use sliding window attention (Child et al., 2019; Beltagy et al., 2020) they can handle generative and embedding inputs of arbitrary length.

**b) Efficiency:** Generative and embedding models are commonly used together to make up for each other’s deficiencies (Guu et al., 2020; Lewis et al., 2021a). One such scenario is Retrieval-Augmented Generation (RAG) (Lewis et al., 2021a), where an embedding model is used to retrieve context that is provided to the generative model to answer a user query. This requires passing the user query and the context into both the generative and the embedding model for a total of four forward passes. With

GRITLM, the embedding and generative model are equivalent, allowing us to cache computations and halve the necessary number of forward passes. We find that this can lead to > 60% faster RAG at inference with long documents.

**c) Simplicity:** Currently, API providers such as OpenAI provide separate generative and embedding endpoints. This requires separate load balancing, additional storage, and more complex serving software. A single model that handles both use cases significantly simplifies infrastructure needs.

Compared to generative instruction tuning, the main downside of GRIT is that it requires more finetuning compute due to training with two objectives. However, finetuning is generally cheap compared to pretraining, thus we think the benefits vastly outstrip this problem. Further, when considering training a separate generative and embedding model from scratch (e.g. for RAG), GritLM is generally cheaper when incorporating the pretraining compute, as there is only one pretraining and finetuning for GritLM, not separate ones for both a generative and an embedding model. Thus we recommend practitioners building instruction-following language models to adopt GRIT.

Alongside GRIT, we introduce novel performance improvements for embedding models including the use of bidirectional attention with mean pooling for LLM embeddings and ensuring that in-batch negatives stem from the same dataset rather than any dataset, as well as novelties for generative models including mixing sample- and token-level loss aggregation. We ablate these in detail in [Appendix A](#). We also put forth new ways to reduce memory requirements during training of embedding models, which we elaborate in [Appendix K](#).

## 2 GRIT

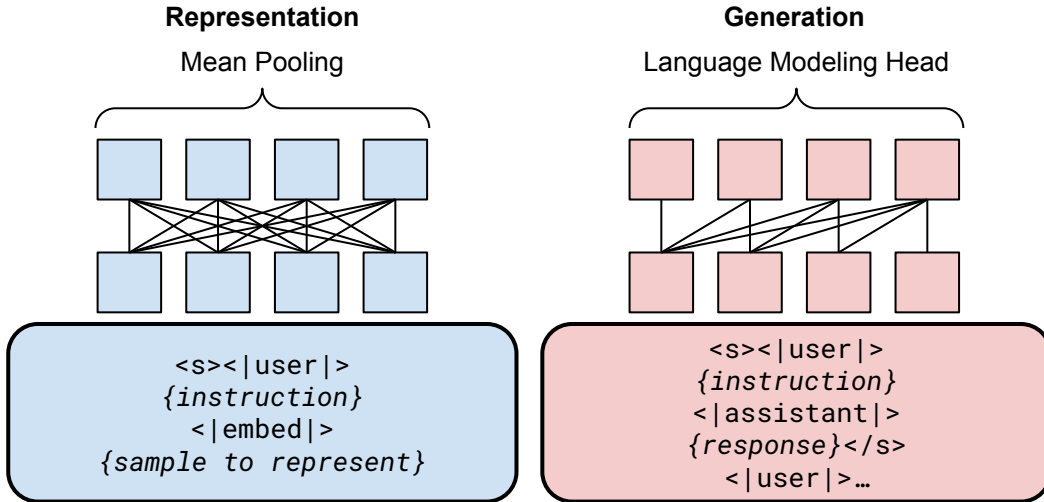


Figure 3: **GRITLM architecture and format.** *Left:* GRITLM uses bidirectional attention over the input for embedding tasks. Mean pooling is applied over the final hidden state to yield the final representation. *Right:* GRITLM uses causal attention over the input for generative tasks. A language modeling head on top of the hidden states predicts the next tokens. The format supports conversations with multiple turns (indicated with “...”).

GRIT unifies representational instruction tuning (Su et al., 2023; Asai et al., 2022; Wang et al., 2024) and generative instruction tuning (Wei et al., 2022; Sanh et al., 2022; Muennighoff et al., 2023d) into a single model. We finetune a pretrained LLM (Brown et al., 2020) with embedding and generative instruction data in a consistent format (Figure 3). For embedding data, we follow prior work and use a contrastive objective with in-batch negatives (Chen et al., 2020; Gao et al., 2022):

$$\mathcal{L}_{\text{Rep}} = -\frac{1}{M} \sum_{i=1}^M \log \frac{\exp(\tau \cdot \sigma(f_{\theta}(q^{(i)}), f_{\theta}(d^{(i)})))}{\sum_{j=1}^M \exp(\tau \cdot \sigma(f_{\theta}(q^{(i)}), f_{\theta}(d^{(j)})))} \quad (1)$$

where  $f$  is GRITLM parametrized by the model  $\theta$ ,  $\tau$  is a temperature hyperparameter and  $\sigma$  corresponds to pooling applied to each output followed by cosine similarity.  $q$  and  $d$  are query and

document samples. As depicted in Figure 3, we use bidirectional attention followed by mean pooling, which corresponds to averaging the hidden states across the sequence length. During pooling, we only average the final hidden states of the input sample, ignoring the instruction and format tokens. However, the instruction and format tokens still influence the final representation through the self-attention mechanism (Vaswani et al., 2023).

To compute the loss on generative data, we use the language modeling objective whereby the model needs to predict the next token (Radford et al., 2018; 2019):

$$\mathcal{L}_{\text{Gen}} = -\frac{1}{N} \sum_{i=1}^N \log P(f_{\theta, \eta}(x^{(i)}) | f_{\theta, \eta}(x^{(<i>i))}) \quad (2)$$

where  $f$  is GRITLM with parameters  $\theta$  and language modeling head  $\eta$ , which is only used for generation.  $x$  are generative training samples. We only compute loss over predicted tokens i.e. “{response}</s>” in Figure 3. A key consideration is how to aggregate the generative loss. Aggregating at the sample level corresponds to giving each sample the same weight within a batch regardless of its token count. Such aggregation is commonly used for instruction tuning, as it can boost performance on discriminative tasks (Muennighoff et al., 2023d). However, Muennighoff et al. (2023d) also show how this in turn can lead to a model biased toward short generations. Meanwhile, aggregation at the token level corresponds to giving each token the same weight, thus samples with many tokens become more important. This leads to a model producing longer generations, which can be important for performance on generative tasks. Especially, human or machine-evaluated generative tasks, such as AlpacaEval (Li et al., 2023b), are known to be biased toward preferring longer generations (Wang et al., 2023). Note that when every sample has the same sequence length such as during pretraining or when the batch size is 1, token and sample level generative loss are equal to each other. One can mix the two to balance their trade-offs, for example doing token level loss across a subset of the batch and then giving each subset the same weight. We explore the trade-offs in our ablations in Appendix A. We sum the objectives with optional loss weights  $\lambda_{\text{Rep}}$  and  $\lambda_{\text{Gen}}$ :

$$\mathcal{L}_{\text{GRIT}} = \lambda_{\text{Rep}} \mathcal{L}_{\text{Rep}} + \lambda_{\text{Gen}} \mathcal{L}_{\text{Gen}} \quad (3)$$

Notably, our formulation supports differing numbers of embedding samples ( $M$ ) and generative samples/tokens ( $N$ ). This allows for significantly increasing the embedding batch size while keeping the generative batch size fixed. A large embedding batch size is often key to well-performing text embedding models (Xiao et al., 2023) at the cost of requiring more compute at each step.

### 3 EXPERIMENTS

In this section, we first outline our experimental setup in §3.1. In §3.2, we discuss and benchmark the embedding and generative performance of our models. In Appendix A, we ablate the settings that led to our final models, including training data, precision, pooling, sequence length, and loss weights.

#### 3.1 SETUP

We finetune our final models from Mistral 7B (Jiang et al., 2023) and Mixtral 8x7B (Jiang et al., 2024) using adaptations of E5 (Wang et al., 2024) and the Tülu 2 data (Iverson et al., 2023). For E5, we adapt it by adding S2ORC (Lo et al., 2020) to increase its scientific data (“E5S”), while for Tülu 2 we filter out their custom prompts that contain answers related to the origin of their model. For GRITLM 7B, we use a batch size of 2048 for embedding data and 256 for generative data and we train the model for a total of 1253 steps corresponding to one epoch on the generative data and 1.36 epochs on the embedding data. For GRITLM 8x7B, the embedding batch size is 256 due to compute limitations. We use several strategies to reduce the memory required during training including a novel technique to split the embedding triplet into separate forward and backward passes detailed in Appendix K. Other hyperparameters are detailed in the ablation experiments in Appendix A and Appendix L. For embedding performance we evaluate using the 56 main datasets from MTEB (Muennighoff et al., 2023c). For generative performance, we largely follow the evaluation setup of Iverson et al. (2023) except that we use the HumanEvalSynthesize (Muennighoff et al., 2023a) variant of HumanEval, as it is more adequate for instruction-following models. We explain each task in detail in Appendix H.

Table 1: **Embedding performance of GRITLM and others.** We indicate parameter counts where available (B=billions). See Appendix H for task, metric, and dataset details. Appendix J contains per-dataset results of GRITLM models. LLMs not finetuned for embedding (Llama 2 70B, Mistral 7B (Instruct), GPT-J 6B, Gen.-only) are evaluated with weighted-mean pooling (Muennighoff, 2022). ♥Results from the MTEB leaderboard (<https://hf.co/spaces/mteb/leaderboard>)

Task (→)	CLF	Clust.	PairCLF	Rerank	Retrieval	STS	Summ.	Avg.
Metric (→)	Acc.	V-Meas.	AP	MAP	nDCG	Spear.	Spear.	
Dataset # (→)	12	11	3	4	15	10	1	56
Proprietary models ♥								
OpenAI v3	75.5	49.0	85.7	59.2	55.4	81.7	29.9	64.6
Other Open Models ♥								
Llama 2 70B	60.4	29.0	47.1	38.5	9.0	49.1	26.1	35.6
Mistral 7B	63.5	34.6	53.5	43.2	13.2	57.4	19.7	40.5
Mistral 7B Instruct	67.1	34.6	59.6	44.8	16.3	63.4	25.9	43.7
GPT-J 6B	66.2	39.0	60.6	48.9	19.8	60.9	26.3	45.2
SGPT BE 5.8B	68.1	40.3	82.0	56.6	50.3	78.1	31.5	58.9
Instructor XL 1.5B	73.1	44.7	86.6	57.3	49.3	83.1	<b>32.3</b>	61.8
BGE Large 0.34B	76.0	46.1	87.1	60.0	54.3	83.1	<u>31.6</u>	64.2
E5 Mistral 7B	78.5	50.3	<b>88.3</b>	60.2	56.9	<b>84.6</b>	31.4	66.6
<b>GRITLM</b>								
Gen.-only 7B	65.4	32.7	54.2	43.0	13.7	60.2	21.1	41.2
Emb.-only 7B	78.8	<b>51.1</b>	87.1	<b>60.7</b>	<b>57.5</b>	83.8	30.2	<b>66.8</b>
GRITLM 7B	<b>79.5</b>	<u>50.6</u>	<u>87.2</u>	<u>60.5</u>	<u>57.4</u>	83.4	30.4	<b>66.8</b>
GRITLM 8x7B	78.5	50.1	85.0	59.8	55.1	83.3	29.8	65.7

Table 2: **Generative performance of GRITLM and others.** We indicate parameter counts where available (B=billions). See Appendix H for dataset, setup, and metric details. ♥Results from Ivison et al. (2023) except for numbers marked with ♦ which are from Touvron et al. (2023) and † which are from us. For models that cannot be easily used as chat models, we set Alpaca to 0.

Dataset (→)	MMLU	GSM8K	BBH	TyDi QA	HumanEval	Alpaca	Avg.
Setup (→)	0 FS	8 FS, CoT	3 FS, CoT	1 FS, GP	0 FS	0 FS, 1.0	
Metric (→)	EM	EM	EM	F1	pass@1	% Win	
Proprietary models ♥							
GPT-4-0613	81.4	95.0	89.1	65.2	86.6 <sup>†</sup>	91.2	84.8
Other Open Models ♥							
GPT-J 6B	27.7	2.5	30.2	9.4	9.8	0.0	13.3
SGPT BE 5.8B	24.4	1.0	0.0	22.8	0.0	0.0	8.0
Zephyr 7B $\beta$	58.6	28.0	44.9	23.7	28.5	85.8	44.9
Llama 2 70B	64.5	55.5	66.0	<b>62.6</b>	29.9 <sup>♦</sup>	0.0	46.4
Llama 2 Chat 13B	53.2	9.0	40.3	32.1	19.6 <sup>†</sup>	91.4	40.9
Llama 2 Chat 70B	60.9	59.0	49.0	44.4	34.3 <sup>†</sup>	<u>94.5</u>	57.0
Tülu 2 7B	50.4	34.0	48.5	46.4	24.5 <sup>†</sup>	73.9	46.3
Tülu 2 13B	55.4	46.0	49.5	53.2	31.4	78.9	52.4
Tülu 2 70B	<u>67.3</u>	<b>73.0</b>	<u>68.4</u>	53.6	41.6	86.6	<u>65.1</u>
Mistral 7B	60.1	44.5	55.6	55.8	30.5	0.0	41.1
Mistral 7B Instruct	53.0	36.0	38.5	27.8	34.0	75.3	44.1
Mixtral 8x7B Instruct	<b>68.4</b>	<u>65.0</u>	55.9	24.3	<b>53.5</b>	<b>94.8</b>	60.3
<b>GRITLM</b>							
Emb.-only 7B	23.5	1.0	0.0	21.0	0.0	0.0	7.6
Gen.-only 7B	57.5	52.0	55.4	56.6	34.5	75.4	55.2
GRITLM 7B	57.6	57.5	54.8	55.4	32.8	74.8	55.5
GRITLM 8x7B	66.7	61.5	<b>70.2</b>	<u>58.2</u>	<u>53.4</u>	84.0	<b>65.7</b>

### 3.2 MAIN RESULTS

**GRIT leads to a strong embedding and generative model** We benchmark GRITLM 7B, GRITLM 8x7B and generative- and embedding-only variants with other models in Table 1 and Table 2. We find that GRITLM 7B outperforms various prior open models on the Massive Text Embedding Benchmark (Muennighoff et al., 2023c) while still outperforming a range of generative models up to its size of 7 billion parameters. For our comparisons we focus on models that are similar to GritLM (e.g. E5 Mistral (Wang et al., 2024) uses the same base model, Instructor (Su et al., 2023) uses a similar dataset, etc.), but we note that there have been various recent embedding and generative models with stronger performance, such as Llama3 (Dubey et al., 2024), NV-Embed (Lee et al., 2024) and others (Li et al., 2024; 2023c; Meng et al., 2024; Chen et al., 2024b; Kim et al., 2024). However, GRIT models are the only ones that can handle both embedding and generation at strong performance (Figure 1). For example, using Llama 70B (Touvron et al., 2023) for embedding leads to a score of only 35.6 on MTEB as depicted in Table 1. GRITLM almost doubles that performance on MTEB, while still outperforming Llama 70B on generative tasks by more than 20% (Table 2). For GRITLM 8x7B, the embedding performance slightly decreases from GRITLM 7B, which is likely because we had to decrease its embedding batch size from 2048 for GRITLM 7B to only 256 for GRITLM 8x7B due to compute limitations (§3.1). We also train embedding-only and generative-only variants of GRITLM that only use representational or generative instruction tuning but are otherwise equivalent. Benchmarking the embedding-only variant (or models like SGPT BE 5.8B (Muennighoff, 2022)) on generative tasks in Table 2 by simply re-adding the language modeling head that was dropped during embedding finetuning leads to around random performance (25.0 is the random baseline on MMLU). Similarly, benchmarking the embedding performance of the generative-only model only leads to a score of 41.2 in Table 1. Thus, joint optimization via the GRIT approach is critical to achieve strong performance for both embedding and generation. We note, however, that with 7 billion parameters GRITLM 7B is significantly more costly to run than many other embedding models in Table 1, such as BGE Large with only 335 million parameters (Xiao et al., 2023). In addition, GRITLM 7B produces representations of 4096 dimensions, which require  $4\times$  more storage than the 1024-dimensional embeddings of BGE Large.

**GRITLM matches embedding-only and generative-only variants** We find that unifying the two objectives via GRITLM matches both the generative-only and the embedding-only variants. This is similar to observations made for visual models (Yu et al., 2022). However, while GRITLM is trained for the same number of steps as the embedding- and generative-only models, it needs more compute per training step as it does a forward and backward pass on both embedding and generative data.

## 4 RERANKING WITH GRIT

Table 3: **Reranking (Rerank) using GRITLM as both Bi- and Cross-Encoder.**

MTEB DS ( $\downarrow$ )	No Rerank	Rerank top 10
ArguAna	63.24	<b>64.39</b>
ClimateFEVER	30.91	<b>31.85</b>
CQADupstack	49.42	<b>50.05</b>
DBPedia	46.60	<b>47.82</b>
FiQA2018	59.95	<b>60.39</b>
FEVER	82.74	<b>82.85</b>
HotpotQA	79.40	<b>80.46</b>
NFCorpus	40.89	<b>41.23</b>
NQ	70.30	<b>71.49</b>
MSMARCO	41.96	<b>42.47</b>
QuoraRetrieval	<b>89.47</b>	88.67
SCIDOCS	24.41	<b>24.54</b>
SciFact	79.17	<b>79.28</b>
TRECCOVID	74.80	<b>75.24</b>
Touche2020	27.93	<b>28.41</b>
Average	57.4	<b>57.9</b>

For retrieval tasks, it is common to follow the embedding-based retrieval stage by a reranking stage (Nogueira & Cho, 2020). In the reranking stage, for each query, the top- $k$  chosen documents are reranked based on a usually more expensive but more performant method. For LLMs, prior work has shown that this can be done by passing each of the  $k$  documents together with the query to the model and scoring the pair with log probabilities (Muennighoff, 2022). Note that this scales quadratically with the number of documents and queries and is thus usually too expensive for the first stage (“Cross-Encoder”). Meanwhile, using embeddings for the first stage is much cheaper as it only requires passing each query and each document once and thus scales linearly (“Bi-Encoder”). More recent work relies on instructions to use LLMs for reranking (Sun et al., 2023; Ma et al., 2023b; Pradeep et al., 2023a;b). While prior work uses separate models for the embedding and rerank-

ing stages, GRITLM can be used for both stages due to its unified capabilities. In Table 3, we display the embedding performance of GRITLM 7B when additionally allowing it to rerank the top 10 documents selected via its embedding capabilities for each query. For reranking, we use the model’s generative capabilities following the permutation generation approach from Sun et al. (2023) and reusing their prompt. We find that reranking via the generative capabilities of GRITLM 7B allows it to improve on its own embedding performance on almost every retrieval dataset. Increasing the top- $k$  documents beyond ten is likely to further improve results, however, at the cost of more compute (Muennighoff, 2022).

## 5 RAG WITH GRIT

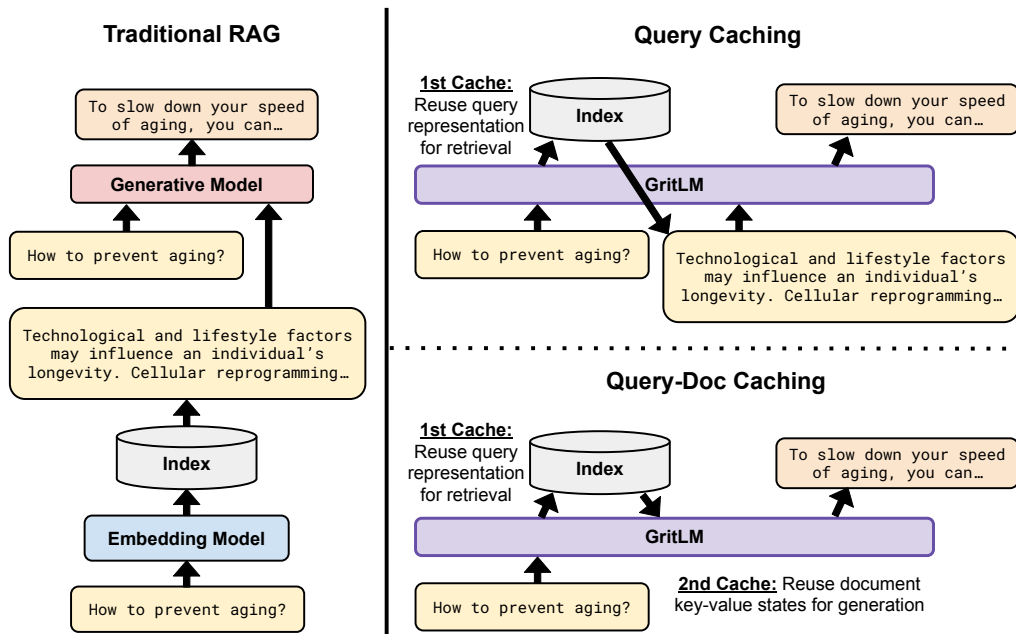


Figure 4: **RAG with GRIT.** *Left:* Traditional Retrieval-Augmented Generation (RAG) relies on a separate embedding model and generative model. *Right:* GRITLM simplifies RAG as it handles both embedding and generation. Query Caching removes the duplicate forward pass of the query by reusing its representation. Query-Doc Caching also removes the forward pass on the document during inference, as the cached index also stores the document key-value states.

**Method** By unifying embedding and generation, GRITLM simplifies Retrieval-Augmented Generation (RAG). Figure 4 displays how caching can reduce forward passes. Specifically, we introduce:   
**(a) Query Caching:** In traditional RAG, the query needs to be passed both through the embedding model and later through the generative model. In Query Caching, we cache the key-value states from the embedding forward pass and reuse them for the generative pass, exploiting the property that both are the same model: GRITLM. Thus, we save compute equivalent to one forward pass of the query. Equivalently, we can also perform the generative forward pass over the query first and use its representation to retrieve the document on the fly (depicted in Figure 4). To make the generations with Query Caching completely equivalent to RAG, we place the query at the beginning of the prompt such that it only attends to itself through causal attention.   
**(b) Doc Caching:** Here we cache the documents,  $D$ . When the index is created, we also save the key-value states of every document and add them to the index. Thus, the index consists of the document embeddings and key-value states. Note that the computational cost of creating the index remains the same as the key-value states have to be computed even if only embeddings are desired. At inference, we still retrieve based on embedding similarity but the index returns the key-value states instead of the text passage. These key-value states are then provided to the model to avoid having to

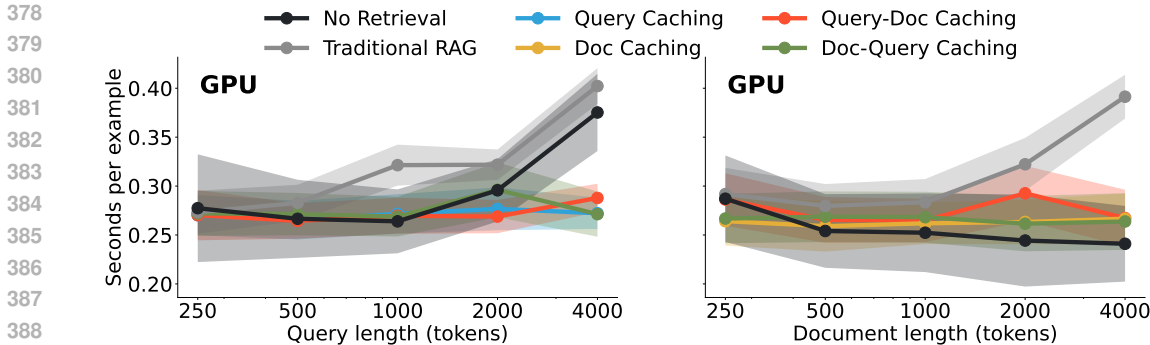


Figure 5: **Inference latency of RAG with GRITLM 7B.** When benchmarking scaling query length (left), document length is fixed at 1, whereas query length is fixed at 1 when scaling document length (right). In addition to the query/doc lengths, the formatting and prompt take up around 40 tokens. We visualize the standard deviation across 100 runs as the shaded area. For each approach, we generate 16 tokens. See Figure 6 for CPU latency.

recompute them. This effectively saves a forward pass for every in-context document at inference. However, this method increases the necessary storage. While the text passages no longer need to be stored, the key-value states now need to be stored and they usually require more storage depending on the model. We note that Document Caching also works for models other than GRITLM. However, for such models, one needs to pass all documents through the generation model ahead of time, thus increasing the cost of creating the index. To maintain equivalence with RAG, the document should be at the beginning of the prompt for Document Caching (opposite of Query Caching).

**(b) Query-Doc Caching / Doc-Query Caching:** We can also combine Query Caching and Doc Caching to save even more inference costs. However, combining them inevitably leads to discrepancies compared to RAG, as in traditional RAG either the query or the document is conditioned on the other one. Meanwhile, if both are cached then they are not conditioned on one another via the self-attention mechanism. We refer to Query-Doc Caching if the query is followed by the document in the prompt and to Doc-Query Caching if the document comes first.

**Setup** We benchmark the caching variants on Natural Questions (Kwiatkowski et al., 2019) using 2,681,468 documents from BEIR NQ (Thakur et al., 2021) as our index. We score models by checking if any correct answer is anywhere in the generation (“match”). Prior work often checks if the generation exactly matches the answer (“exact match”) (Izacard et al., 2022). However, due to the chat data our model answers in few sentences, thus exact match fails to credit many correct answers. In the first 20 samples of the “No RAG” baseline, “exact match” leads to 4 false negatives that “match” credits correctly without any false positives. We do not use instructions for embedding here, only the format in Figure 3.

**Performance** As depicted in Table 4, RAG performs better than the “No RAG” baseline where the model is not provided any context. This validates that despite its small size compared to prior work (Lin et al., 2023), our index is still valuable. While Query and Doc Caching can theoretically lead to the exact same performance as RAG, we experience differences for two reasons: **1) Attention:** Our model is trained to embed with bidirectional attention (§2) and thus we use bidirectional attention when embedding query or document. Meanwhile, the generative model expects causal key-value states. In the Query-Doc/Doc-Query setup, there is an additional mismatch in either the documents or the queries not having attended to the other one, as both need to be embedded and cached separately. **2) Formatting:** The query is formatted in the embedding format as depicted in Figure 3, which the model has never seen during generative training. This could further lead to a performance drop. Due to 1) and 2), Query Caching leads to a performance drop compared to traditional RAG. However, the Query Caching performance of 25.46 is still better than not using RAG, thus it comes down to a speed-performance trade-off. Formatting the RAG baseline using the embedding format (Figure 3) reduces its score from 30.50 to 29.36 (not depicted), thus the additional four-point discrepancy of Query Caching and the majority of the damage is because of the attention issue. Meanwhile, Doc Caching slightly improves performance resulting in the best match score among all methods



Table 4: **RAG benchmarking on Natural Questions with GRITLM 7B**. For RAG, the retrieved context is simply placed in the context of the language model in contrast to our caching alternatives (Figure 4). CPU and GPU latencies are measured on an “Intel(R) Xeon(R) Platinum 8481C CPU @ 2.70GHz” and one “NVIDIA H100 80GB HBM3”, respectively. Sample A has a query of 1 token and a document of 4000 tokens, and sample B is the inverse. For each approach, we generate 16 tokens. Storage consists of the index and passages, except for Doc Caching variants where it is the index and key-value states. The index is stored in float32, while key-value states are stored in bfloat16. See Appendix F for experiments on TriviaQA and MMLU.

	Match (0-shot, ↑)	CPU Latency (s, ↓)		GPU Latency (s, ↓)		Storage (↓)
		Sample A	Sample B	Sample A	Sample B	
No RAG	21.00	4.3 ± 0.36	13.69 ± 1.0	0.24 ± 0.04	0.38 ± 0.04	<b>0GB</b>
<i>Query then document prompt</i>						
RAG	<u>30.50</u>	11.64 ± 0.74	14.88 ± 0.87	0.39 ± 0.02	0.40 ± 0.02	<u>43GB</u>
Query Caching	25.46	18.30 ± 0.76	6.87 ± 0.89	0.44 ± 0.03	<b>0.27 ± 0.02</b>	<u>43GB</u>
Query-Doc Caching	21.63	<b>5.12 ± 0.23</b>	<u>6.62 ± 0.97</u>	<u>0.27 ± 0.03</u>	0.29 ± 0.01	<u>30TB</u>
<i>Document then query prompt</i>						
RAG	30.47	14.18 ± 1.01	15.33 ± 0.87	0.39 ± 0.01	0.4 ± 0.01	<u>43GB</u>
Doc Caching	<b>33.38</b>	5.25 ± 0.34	23.23 ± 1.05	<u>0.27 ± 0.03</u>	0.45 ± 0.02	<u>30TB</u>
Doc-Query Caching	18.39	<u>5.23 ± 0.37</u>	<b>6.41 ± 0.96</b>	<b>0.26 ± 0.03</b>	<b>0.27 ± 0.02</b>	<u>30TB</u>

considered. This is possibly because, unlike the query, the document does not need to be as thoroughly understood, and skimming it may suffice. Thus, the slightly corrupted key-value states do not result in a performance drop. Query-Doc and Doc-Query Caching only perform near the “No RAG” baseline in our experiments, which may limit their usefulness in practice. This is likely caused by the additional attention mismatch that they introduce. This issue as well as the formatting issue could likely be solved by an additional RAG finetuning stage on top of GRITLM, which we leave to future work.

**Latency** Caching is much faster than RAG on both CPUs and GPUs, especially for long sequences (Figure 5). In Table 4, we display that for 4000 tokens, Query Caching is 54% and 33% faster on CPUs and GPUs, respectively (Sample B). For Doc Caching it is 63% and 31% (Sample A). If going beyond 4000 tokens the speed-ups will be even larger. However, for the opposite samples in Table 4 speed remains around the same. This is because while for Sample A, Doc Caching caches 4000 tokens, for Sample B it caches only 1 token, which does not provide any speed-up. Thus, Doc Caching should be used when documents are expected to be very long, while Query Caching should be used when queries are expected to be very long. In a production setting, a simple input length check could switch from one caching mode to the other. As is the case in Table 4, caching can match or even be faster than not using retrieval at all (“No RAG”). This could be due to the embedding forward pass not using the language modeling head. For Query Caching, the language modeling head is only used for the tokens that are generated, while for “RAG” and “No RAG” it is used for the entire input. The matrix multiplication with the language modeling head is computationally expensive due to its high dimensionality, which could cause the slower speed of the no retrieval baseline. Query-Doc Caching and Doc-Query Caching cache both documents and queries and thus lead to major speed-ups for both Sample A and Sample B in Table 4. Overall, speed-ups are larger on CPUs, as GPUs can process the entire sequence in parallel, thus the advantage of caching parts of it is smaller. We also note that our RAG baseline uses our 7B parameter model for both the embedding and generative model but without caching. In practice, it is often common to have an embedding model that is much smaller and cheaper than the generative model. Nonetheless, as caching with GRITLM-7B approaches the No RAG latency in Table 4, we still expect it to be faster than setups with smaller embedding models for long sequences. In addition, it would lead to significantly better performance in that case due to the state-of-the-art retrieval performance of GRITLM.

**Storage** In most RAG setups the embeddings of all documents are precomputed and stored to be later used at inference. This is referred to as the index. In traditional RAG, the documents themselves still need to be stored, as the index is only used to find the document ID, which is then used to fetch the document text and pass it to the generative model. For Doc Caching variants documents no longer need to be stored, however, the key-value states need to be stored. The key-value states take

up a lot of storage, as they consist of two tensors of shape (batch size, number of heads, sequence length, dimension per head) for each batch. For our 2,681,468 documents and the 7-billion parameter GRITLM model, this leads to 30TB of key-value states. However, unlike the index, the key-value states can be fully offloaded to disk and do not need to be kept in memory. Once the document ID has been determined via the index, the corresponding key-value state can be simply loaded from disk. For a single sample, this corresponds to loading 12.5MB of key-value states into memory.

## 6 RELATED WORK

The story of text embedding and text generation has been a story of unification.

**Embedding Models** used to focus on word representations (Pennington et al., 2014; Mikolov et al., 2013) that struggled generalizing to entire sentences or passages (Conneau & Kiela, 2018). Inference (Conneau et al., 2018), SBERT (Reimers & Gurevych, 2019) and similar models (Ni et al., 2021b;a) emerged that handle both the embedding of words and sentences at good quality by considering context when present. However, for strong performance, they require separate models for symmetric and asymmetric tasks (Muennighoff et al., 2023c; Neelakantan et al., 2022). Symmetric embedding tasks are ones where the query and document are expected to come from the same distribution, such as STS. Meanwhile, for asymmetric tasks, they come from different distributions and as such could have very different sequence lengths like in retrieval. For example, the MTEB benchmark (Muennighoff et al., 2023c) revealed that SentT5 (Ni et al., 2021b) only performs well at symmetric tasks, while GTR (Ni et al., 2021a) only at asymmetric tasks despite both using T5 (Raffel et al., 2023) as their base model. Recent embedding models have been able to unify symmetric and asymmetric tasks into a single model by differentiating them in the prompt (Xiao et al., 2023; Wang et al., 2022a). Further, including detailed instructions in the prompt has allowed unifying practically any embedding task into a single model (Su et al., 2023).

**Generative Models** used to be tailored to a single task, such as translation (Sutskever et al., 2014) or question answering (Yin et al., 2016). McCann et al. (2018) cast multiple generative tasks as question answering to unify them within a single model, however, performance was still limited and it did not generalize to arbitrary tasks. Large-scale self-supervised pretraining has enabled the use of a single large language model (LLM) for practically any generative task (Brown et al., 2020; Chowdhery et al., 2022; Rae et al., 2022; BigScience Workshop et al., 2023; Scao et al., 2022; Groeneveld et al., 2024; Li et al., 2023a). However, using an LLM without careful prompting often leads to poor performance (Rubin et al., 2022; Min et al., 2022b). Finetuning LLMs on instructions has emerged as a method to significantly ease the usage of the models to apply them to any generative task with strong results (Wei et al., 2022; Sanh et al., 2022; Min et al., 2022a; Wang et al., 2022c; Mishra et al., 2022; Iyer et al., 2023; Üstün et al., 2024; Singh et al., 2024; Zhou et al., 2023).

The two streams of embedding and generative models have each been unified into a single model that handles any task within its stream. Unifying the two streams into a single model that handles any task both for embedding and generation is the natural next step toward a general multi-task model.

## 7 CONCLUSION

We present GRIT to unify text embedding and generation, and thus all text-based language problems, into one model: GRITLM. GRITLM 7B performs strongly on the Massive Text Embedding Benchmark, while simultaneously possessing generative capabilities that exceed some larger models. Notably, its performance matches otherwise equivalent embedding-only and generative-only variants allowing us to unify them at no performance loss. We show that GRIT simplifies the field using the examples of reranking and RAG. For reranking, we are able to improve retrieval performance by around 10% by reusing GRITLM as reranker instead of having to rely on a separate model. For RAG, we unify the retriever and reader into a single model, GRITLM, speeding up inference by >60% for long texts at no performance loss via GRIT Doc Caching. We believe GRIT paves the way for a paradigm shift in language modeling, where embedding and generation seamlessly coexist in a single model. As such, we highlight the various limitations of this work and point the community to potential future research in [Appendix Q](#).

## REFERENCES

- 540  
541  
542 Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos,  
543 Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report, 2023.
- 544  
545 Sunil Arya, David M Mount, Nathan S Netanyahu, Ruth Silverman, and Angela Y  
546 Wu. An optimal algorithm for approximate nearest neighbor searching fixed dimen-  
547 sions, 1998. URL [https://graphics.stanford.edu/courses/cs468-06-fall/  
548 Papers/03%20AMNSW%20-%20JACM.pdf](https://graphics.stanford.edu/courses/cs468-06-fall/Papers/03%20AMNSW%20-%20JACM.pdf).
- 549  
550 Akari Asai, Timo Schick, Patrick Lewis, Xilun Chen, Gautier Izacard, Sebastian Riedel, Hannaneh  
551 Hajishirzi, and Wen tau Yih. Task-aware retrieval with instructions, 2022.
- 552  
553 Akari Asai, Sewon Min, Zexuan Zhong, and Danqi Chen. Retrieval-based language models and  
554 applications, 2023a. URL <https://aclanthology.org/2023.acl-tutorials.6/>.
- 555  
556 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to  
557 retrieve, generate, and critique through self-reflection, 2023b.
- 558  
559 Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016.
- 560  
561 Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Ma-  
562 jumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, Mir Rosenberg, Xia Song, Alina Stoica,  
563 Saurabh Tiwary, and Tong Wang. Ms marco: A human generated machine reading comprehension  
564 dataset, 2018.
- 565  
566 Parishad BehnamGhader, Vaibhav Adlakha, Marius Mosbach, Dzmitry Bahdanau, Nicolas Chapados,  
567 and Siva Reddy. Llm2vec: Large language models are secretly powerful text encoders. *arXiv  
568 preprint arXiv:2404.05961*, 2024.
- 569  
570 Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer,  
571 2020.
- 572  
573 Loubna Ben Allal, Niklas Muennighoff, Logesh Kumar Umapathi, Ben Lipkin, and Leandro von  
574 Werra. A framework for the evaluation of code generation models, 2022. URL [https://  
575 github.com/bigcode-project/bigcode-evaluation-harness](https://github.com/bigcode-project/bigcode-evaluation-harness).
- 576  
577 Teven Le Scao BigScience Workshop, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić,  
578 Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé,  
579 Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammana-  
580 manchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, et al. Bloom: A 176b-parameter  
581 open-access multilingual language model, 2023.
- 582  
583 Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican,  
584 George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. Improving  
585 language models by retrieving from trillions of tokens, 2022.
- 586  
587 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,  
588 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
589 few-shot learners, 2020.
- 590  
591 Isabel Cachola, Kyle Lo, Arman Cohan, and Daniel S. Weld. Tldr: Extreme summarization of  
592 scientific documents, 2020.
- 593  
594 Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. Bge m3-embedding:  
595 Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge  
596 distillation, 2024a.
- 597  
598 Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. Bge m3-embedding:  
599 Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge  
600 distillation, 2024b. URL <https://arxiv.org/abs/2402.03216>.
- 601  
602 Lingjiao Chen, Matei Zaharia, and James Zou. How is chatgpt’s behavior changing over time?, 2023.

- 594 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared  
595 Kaplan, Harri Edwards, Yuri Burda, et al. Evaluating large language models trained on code, 2021.  
596
- 597 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for  
598 contrastive learning of visual representations, 2020.
- 599 Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollar, and  
600 C. Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server, 2015.  
601
- 602 Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse  
603 transformers, 2019.
- 604 Jaemin Cho, Jie Lei, Hao Tan, and Mohit Bansal. Unifying vision-and-language tasks via text  
605 generation, 2021.  
606
- 607 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam  
608 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:  
609 Scaling language modeling with pathways, 2022.
- 610 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li,  
611 Xuezhi Wang, et al. Scaling instruction-finetuned language models, 2022.  
612
- 613 Jonathan H. Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev,  
614 and Jennimaria Palomaki. Tydi qa: A benchmark for information-seeking question answering in  
615 typologically diverse languages, 2020.  
616
- 617 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
618 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John  
619 Schulman. Training verifiers to solve math word problems, 2021.
- 620 Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel S. Weld. Specter: Document-  
621 level representation learning using citation-informed transformers, 2020.  
622
- 623 Alexis Conneau and Douwe Kiela. Senteval: An evaluation toolkit for universal sentence representa-  
624 tions, 2018.
- 625 Alexis Conneau, Douwe Kiela, Holger Schwenk, Loic Barrault, and Antoine Bordes. Supervised  
626 learning of universal sentence representations from natural language inference data, 2018.  
627
- 628 William Coster and David Kauchak. Simple english wikipedia: A new text simplification task, 2011.  
629 URL <https://api.semanticscholar.org/CorpusID:9128245>.
- 630 Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning, 2023.  
631
- 632 Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and  
633 memory-efficient exact attention with io-awareness, 2022.
- 634 DataCanary, hilfialkaff, Meg Risdal Lili Jiang, Nikhil Dandekar, and tomtung. Quora question pairs,  
635 2017. URL <https://kaggle.com/competitions/quora-question-pairs>.  
636
- 637 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep  
638 bidirectional transformers for language understanding, 2019.  
639
- 640 Kaustubh D. Dhole, Varun Gangal, Sebastian Gehrmann, Aadesh Gupta, Zhenhao Li, Saad Ma-  
641 hamood, Abinaya Mahendiran, Simon Mille, Ashish Shrivastava, Samson Tan, et al. Nl-augmenter:  
642 A framework for task-sensitive natural language augmentation, 2022.
- 643 Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong  
644 Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional  
645 conversations, 2023.  
646
- 647 Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvassy, Pierre-Emmanuel  
Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library, 2024.

648 Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. All  
649 nlp tasks are generation tasks: A general pretraining framework, 2021. URL <https://arxiv.org/abs/2103.10360v1>.  
650  
651

652 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
653 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn,  
654 Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston  
655 Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron,  
656 Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris  
657 McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton  
658 Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David  
659 Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes,  
660 Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip  
661 Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme  
662 Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu,  
663 Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov,  
664 Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah,  
665 Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu  
666 Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph  
667 Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani,  
668 Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz  
669 Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence  
670 Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas  
671 Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri,  
672 Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis,  
673 Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov,  
674 Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Celebi, Patrick Alrassy, Pengchuan  
675 Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan,  
676 Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy,  
677 Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit  
678 Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou,  
679 Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia  
680 Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparth, Sheng Shen, Shengye Wan,  
681 Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla,  
682 Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek  
683 Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao,  
684 Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent  
685 Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu,  
686 Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia,  
687 Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen  
688 Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe  
689 Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya  
690 Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex  
691 Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei  
692 Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew  
693 Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley  
694 Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin  
695 Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu,  
696 Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt  
697 Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changan Wang, Changkyu Kim, Chao  
698 Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon  
699 Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide  
700 Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le,  
701 Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily  
Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix  
Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank  
Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern,  
Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid  
Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen

- 702 Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-  
703 Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste  
704 Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul,  
705 Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie,  
706 Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik  
707 Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly  
708 Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhota, Kyle Huang, Lailin Chen,  
709 Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu,  
710 Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria  
711 Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev,  
712 Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle  
713 Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang,  
714 Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam,  
715 Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier,  
716 Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia  
717 Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro  
718 Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani,  
719 Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy,  
720 Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan  
721 Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara  
722 Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh  
723 Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha,  
724 Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe,  
725 Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan  
726 Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury,  
727 Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe  
728 Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi,  
729 Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vitor Albiero, Vlad Ionescu,  
730 Vlad Poenaru, Vlad Tiberiu Mihalescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang,  
731 Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojuan Wu, Xiaolan Wang,  
732 Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang,  
733 Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait,  
734 Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd  
735 of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- 734 Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin,  
735 Percy Liang, and Tatsunori B. Hashimoto. AlpacaFarm: A simulation framework for methods that  
736 learn from human feedback, 2023.
- 737 Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled  
738 alpacaEval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*, 2024.
- 739 Matthew Dunn, Levent Sagun, Mike Higgins, V. Ugur G’üneş, Volkan Cirik, and Kyunghyun  
740 Cho. Searchqa: A new q&a dataset augmented with context from a search engine, 2017. URL  
741 <http://arxiv.org/abs/1704.05179>.
- 742 Paul-Ambroise Duquenne, Holger Schwenk, and Benoît Sagot. Sonar: Sentence-level multimodal  
743 and language-agnostic representations, 2023.
- 744 Hady ElSahar, Pavlos Vougiouklis, Arslan Remaci, Christophe Gravier, Jonathon S. Hare, Frédérique  
745 Laforest, and Elena Paslaru Bontas Simperl. T-rex: A large scale alignment of natural language with  
746 knowledge base triples, 2018. URL [https://api.semanticscholar.org/CorpusID:  
747 4612975](https://api.semanticscholar.org/CorpusID:4612975).
- 748 Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model  
749 alignment as prospect theoretic optimization, 2024.
- 750 Alexander R. Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and  
751 Dragomir Radev. Summeval: Re-evaluating summarization evaluation, 2021.

- 756 Anthony Fader, Luke Zettlemoyer, and Oren Etzioni. Open question answering over curated  
757 and extracted knowledge bases, 2014. URL [https://api.semanticscholar.org/  
758 CorpusID:207214527](https://api.semanticscholar.org/CorpusID:207214527).
- 759 Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. Eli5:  
760 Long form question answering, 2019.
- 762 Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. Language-agnostic  
763 bert sentence embedding, 2022.
- 764 Katja Filippova and Yasemin Altun. Overcoming the lack of parallel data in sentence compression,  
765 2013. URL <https://api.semanticscholar.org/CorpusID:9751546>.
- 767 Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Gold-  
768 ing, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, Jason Phang, Laria Reynolds, Eric Tang,  
769 Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model  
770 evaluation, September 2021a. URL <https://doi.org/10.5281/zenodo.5371628>.
- 771 Luyu Gao, Yunyi Zhang, Jiawei Han, and Jamie Callan. Scaling deep contrastive learning batch size  
772 under memory limited setup, 2021b.
- 774 Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence  
775 embeddings, 2022.
- 776 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Qianyu  
777 Guo, Meng Wang, and Haofen Wang. Retrieval-augmented generation for large language models:  
778 A survey, 2024.
- 779 David Graff, Junbo Kong, Ke Chen, and Kazuaki Maeda. English gigaword, 2003. URL [https://  
781 catalog.ldc.upenn.edu/LDC2011T07](https://catalog.ldc.upenn.edu/LDC2011T07).
- 782 Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafford,  
783 Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson,  
784 Russell Authur, Khyathi Raghavi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu,  
785 Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik,  
786 Crystal Nam, Matthew E. Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk,  
787 Saurabh Shah, Will Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep  
788 Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Sol-  
789 daini, Noah A. Smith, and Hannaneh Hajishirzi. Olmo: Accelerating the science of language  
790 models, 2024.
- 791 Mansi Gupta, Nitish Kulkarni, Raghuveer Chanda, Anirudha Rayasam, and Zachary C Lipton.  
792 Amazonqa: A review-based question answering task, 2019.
- 793 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. Realm: Retrieval-  
794 augmented language model pre-training, 2020.
- 796 Michael Günther, Louis Milliken, Jonathan Geuter, Georgios Mastrapas, Bo Wang, and Han Xiao.  
797 Jina embeddings: A novel set of high-performance sentence embedding models, 2023.
- 798 Michael Günther, Jackmin Ong, Isabelle Mohr, Alaeddine Abdessalem, Tanguy Abel, Moham-  
799 mad Kalim Akram, Susana Guzman, Georgios Mastrapas, Saba Sturua, Bo Wang, Maximilian  
800 Werk, Nan Wang, and Han Xiao. Jina embeddings 2: 8192-token general-purpose text embeddings  
801 for long documents, 2024.
- 802 Felix Hamborg, Norman Meuschke, Corinna Breitingner, and Bela Gipp. news-please - a generic news  
803 crawler and extractor, 2017. URL [https://api.semanticscholar.org/CorpusID:  
804 5830937](https://api.semanticscholar.org/CorpusID:5830937).
- 806 Dan Hendrycks, Nicholas Carlini, John Schulman, and Jacob Steinhardt. Unsolved problems in ml  
807 safety, 2022.
- 808 Christopher Hidey and Kathy McKeown. Identifying causal relations using parallel Wikipedia  
809 articles, August 2016. URL <https://aclanthology.org/P16-1135>.

- 810 Md. Zakir Hossain, Ferdous Sohel, Mohd Fairuz Shiratuddin, and Hamid Laga. A comprehensive  
811 survey of deep learning for image captioning, 2018.  
812
- 813 Jui-Ting Huang, Ashish Sharma, Shuying Sun, Li Xia, David Zhang, Philip Pronin, Janani Padman-  
814 abhan, Giuseppe Ottaviano, and Linjun Yang. Embedding-based retrieval in facebook search, 2020.  
815 URL <https://arxiv.org/abs/2006.11632>.
- 816 Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep Dasigi,  
817 Joel Jang, David Wadden, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. Camels in a  
818 changing climate: Enhancing lm adaptation with tulu 2, 2023.  
819
- 820 Srinivasan Iyer, Xi Victoria Lin, Ramakanth Pasunuru, Todor Mihaylov, Daniel Simig, Ping Yu, Kurt  
821 Shuster, Tianlu Wang, Qing Liu, Punit Singh Koura, Xian Li, Brian O’Horo, Gabriel Pereyra,  
822 Jeff Wang, Christopher Dewan, Asli Celikyilmaz, Luke Zettlemoyer, and Ves Stoyanov. Opt-impl:  
823 Scaling language model instruction meta learning through the lens of generalization, 2023.
- 824 Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane  
825 Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Atlas: Few-shot learning with  
826 retrieval augmented language models, 2022.  
827
- 828 Andrew Jaegle, Felix Gimeno, Andrew Brock, Andrew Zisserman, Oriol Vinyals, and Joao Carreira.  
829 Perceiver: General perception with iterative attention, 2021.
- 830 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
831 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,  
832 L lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas  
833 Wang, Timoth e Lacroix, and William El Sayed. Mistral 7b, 2023.
- 834 Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris  
835 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand,  
836 Gianna Lengyel, Guillaume Bour, Guillaume Lample, L lio Renard Lavaud, Lucile Saulnier, Marie-  
837 Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le  
838 Scao, Th ophile Gervet, Thibaut Lavril, Thomas Wang, Timoth e Lacroix, and William El Sayed.  
839 Mixtral of experts, 2024.  
840
- 841 Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W. Cohen, and Xinghua Lu. Pubmedqa: A  
842 dataset for biomedical research question answering, 2019.
- 843 Jeff Johnson, Matthijs Douze, and Herv  J gou. Billion-scale similarity search with gpus, 2017.  
844
- 845 Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly  
846 supervised challenge dataset for reading comprehension, 2017. URL [https://arxiv.org/  
847 abs/1705.03551](https://arxiv.org/abs/1705.03551).
- 848 Lukasz Kaiser, Aidan N. Gomez, Noam Shazeer, Ashish Vaswani, Niki Parmar, Llion Jones, and  
849 Jakob Uszkoreit. One model to learn them all, 2017.  
850
- 851 Uday Kamath, John Liu, and James Whitaker. Deep learning for nlp and speech recognition, 2019.  
852 URL <https://link.springer.com/book/10.1007/978-3-030-14596-5>.
- 853 Vladimir Karpukhin, Barlas O uz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi  
854 Chen, and Wen tau Yih. Dense passage retrieval for open-domain question answering, 2020.  
855
- 856 Phillip Keung, Yichao Lu, Gy rgy Szarvas, and Noah A. Smith. The multilingual amazon reviews  
857 corpus, 2020.
- 858 Daniel Khashabi, Amos Ng, Tushar Khot, Ashish Sabharwal, Hannaneh Hajishirzi, and Chris  
859 Callison-Burch. Gooaq: Open question answering with diverse answer types, 2021.  
860
- 861 Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Casey A  
862 Fitzpatrick, Peter Bull, Greg Lipstein, Tony Nelli, Ron Zhu, et al. The hateful memes challenge:  
863 Competition report, 2021. URL [https://proceedings.mlr.press/v133/kiela21a.  
html](https://proceedings.mlr.press/v133/kiela21a.html).



- 864 Junseong Kim, Seolhwa Lee, Jihoon Kwon, Sangmo Gu, Yejin Kim, Minkyung Cho, Jy yong  
865 Sohn, and Chanyeol Choi. Linq-embed-mistral:elevating text retrieval with improved gpt data  
866 through task-specific control and quality refinement. Linq AI Research Blog, 2024. URL <https://getlinq.com/blog/linq-embed-mistral/>.  
867
- 868 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.  
869
- 870 Mahnaz Koupaee and William Yang Wang. Wikihow: A large scale text summarization dataset, 2018.  
871
- 872 Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris  
873 Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions:  
874 a benchmark for question answering research, 2019. URL <https://aclanthology.org/Q19-1026/>.  
875
- 876 Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens,  
877 Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri,  
878 David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and  
879 Alexander Mattick. Openassistant conversations – democratizing large language model alignment,  
880 2023.  
881
- 882 Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catanzaro,  
883 and Wei Ping. Nv-embed: Improved techniques for training llms as generalist embedding models,  
884 2024. URL <https://arxiv.org/abs/2405.17428>.
- 885 Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. Zero-shot relation extraction via  
886 reading comprehension, 2017.  
887
- 888 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,  
889 Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela.  
890 Retrieval-augmented generation for knowledge-intensive nlp tasks, 2021a.
- 891 Patrick Lewis, Yuxiang Wu, Linqing Liu, Pasquale Minervini, Heinrich Küttler, Aleksandra Piktus,  
892 Pontus Stenetorp, and Sebastian Riedel. Paq: 65 million probably-asked questions and what you  
893 can do with them, 2021b.  
894
- 895 Chaofan Li, MingHao Qin, Shitao Xiao, Jianlyu Chen, Kun Luo, Yingxia Shao, Defu Lian, and  
896 Zheng Liu. Making text embedders few-shot learners, 2024. URL <https://arxiv.org/abs/2409.15700>.  
897
- 898 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou,  
899 Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. Starcoder: may the source be with  
900 you!, 2023a.  
901
- 902 Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy  
903 Liang, and Tatsunori B. Hashimoto. AlpacaEval: An automatic evaluator of instruction-following  
904 models, 2023b. URL [https://github.com/tatsu-lab/alpaca\\_eval](https://github.com/tatsu-lab/alpaca_eval).
- 905 Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. Towards  
906 general text embeddings with multi-stage contrastive learning, 2023c. URL <https://arxiv.org/abs/2308.03281>.  
907
- 908 Xi Victoria Lin, Xilun Chen, Mingda Chen, Weijia Shi, Maria Lomeli, Rich James, Pedro Rodriguez,  
909 Jacob Kahn, Gergely Szilvasy, Mike Lewis, Luke Zettlemoyer, and Scott Yih. Ra-dit: Retrieval-  
910 augmented dual instruction tuning, 2023.  
911
- 912 Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Dan S. Weld. S2orc: The semantic  
913 scholar open research corpus, 2020.  
914
- 915 Shayne Longpre, Robert Mahari, Anthony Chen, Naana Obeng-Marnu, Damien Sileo, William  
916 Brannon, Niklas Muennighoff, Nathan Khazam, Jad Kabbara, Kartik Perisetla, Xinyi Wu, Enrico  
917 Shippole, Kurt Bollacker, Tongshuang Wu, Luis Villa, Sandy Pentland, Deb Roy, and Sara Hooker.  
The data provenance initiative: A large scale audit of dataset licensing & attribution in ai, 2023.

- 918 Hongyin Luo, Yung-Sung Chuang, Yuan Gong, Tianhua Zhang, Yoon Kim, Xixin Wu, Danny Fox,  
919 Helen Meng, and James Glass. Sail: Search-augmented instruction learning, 2023.  
920
- 921 Risto Luukkonen, Ville Komulainen, Jouni Luoma, Anni Eskelinen, Jenna Kanerva, Hanna-Mari  
922 Kupari, Filip Ginter, Veronika Laippala, Niklas Muennighoff, Aleksandra Piktus, Thomas Wang,  
923 Nouamane Tazi, Teven Le Scao, Thomas Wolf, Osma Suominen, Samuli Sairanen, Mikko Merioksa,  
924 Jyrki Heinonen, Aija Vahtola, Samuel Antao, and Sampo Pyysalo. Fingpt: Large generative models  
925 for a small language, 2023.
- 926 Xueguang Ma, Liang Wang, Nan Yang, Furu Wei, and Jimmy Lin. Fine-tuning llama for multi-stage  
927 text retrieval, 2023a.
- 928 Xueguang Ma, Xinyu Zhang, Ronak Pradeep, and Jimmy Lin. Zero-shot listwise document reranking  
929 with a large language model, 2023b.  
930
- 931 Bryan McCann, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. The natural language  
932 decathlon: Multitask learning as question answering, 2018.
- 933 Rui Meng, Ye Liu, Shafiq Rayhan Joty, Caiming Xiong, Yingbo Zhou, and Semih Yavuz. Sfr-  
934 embedding-mistral:enhance text retrieval with transfer learning. Salesforce AI Research Blog, 2024.  
935 URL <https://blog.salesforceairesearch.com/sfr-embedded-mistral/>.
- 936 Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed representa-  
937 tions of words and phrases and their compositionality, 2013.  
938
- 939 Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hannaneh Hajishirzi. Metaicl: Learning to learn in  
940 context, 2022a.
- 941 Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke  
942 Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work?,  
943 2022b.
- 944 Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. Cross-task generalization  
945 via natural language crowdsourcing instructions, 2022.  
946
- 947 John X. Morris, Volodymyr Kuleshov, Vitaly Shmatikov, and Alexander M. Rush. Text embeddings  
948 reveal (almost) as much as text, 2023.
- 949 Niklas Muennighoff. Vilio: State-of-the-art visio-linguistic models applied to hateful memes, 2020.  
950
- 951 Niklas Muennighoff. Sgpt: Gpt sentence embeddings for semantic search, 2022.
- 952 Niklas Muennighoff, Qian Liu, Armel Zebaze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo, Swayam  
953 Singh, Xiangru Tang, Leandro von Werra, and Shayne Longpre. Octopack: Instruction tuning code  
954 large language models, 2023a.
- 955 Niklas Muennighoff, Alexander M. Rush, Boaz Barak, Teven Le Scao, Aleksandra Piktus, Nouamane  
956 Tazi, Sampo Pyysalo, Thomas Wolf, and Colin Raffel. Scaling data-constrained language models,  
957 2023b.  
958
- 959 Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. Mteb: Massive text embedding  
960 benchmark, 2023c.
- 961 Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le  
962 Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir  
963 Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson,  
964 Edward Raff, and Colin Raffel. Crosslingual generalization through multitask finetuning, 2023d.  
965
- 966 Shashi Narayan, Shay B. Cohen, and Mirella Lapata. Don't give me the details, just the summary!  
967 topic-aware convolutional neural networks for extreme summarization, 2018.
- 968 Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming  
969 Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, Johannes Heidecke, Pranav Shyam, Boris  
970 Power, Tyna Eloundou Nekoul, Girish Sastry, Gretchen Krueger, David Schnurr, Felipe Petroski  
971 Such, Kenny Hsu, Madeleine Thompson, Tabarak Khan, Toki Sherbakov, Joanne Jang, Peter  
Welinder, and Lilian Weng. Text and code embeddings by contrastive pre-training, 2022.

- 972 Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernández Ábrego, Ji Ma, Vincent Y. Zhao,  
973 Yi Luan, Keith B. Hall, Ming-Wei Chang, and Yinfei Yang. Large dual encoders are generalizable  
974 retrievers, 2021a.
- 975  
976 Jianmo Ni, Gustavo Hernández Ábrego, Noah Constant, Ji Ma, Keith B. Hall, Daniel Cer, and Yinfei  
977 Yang. Sentence-t5: Scalable sentence encoders from pre-trained text-to-text models, 2021b.
- 978 Rodrigo Nogueira and Kyunghyun Cho. Passage re-ranking with bert, 2020.
- 979  
980 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni  
981 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, et al. Gpt-4 technical report, 2023.
- 982  
983 Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. Medmcqa : A large-scale  
984 multi-subject multi-choice dataset for medical domain question answering, 2022.
- 985  
986 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor  
987 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward  
988 Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang,  
989 Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning  
990 library, 2019.
- 991  
992 Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word  
993 representation, 2014. URL <https://aclanthology.org/D14-1162/>.
- 994  
995 Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James  
996 Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, Vassilis Plachouras, Tim Rocktäschel,  
997 and Sebastian Riedel. Kilt: a benchmark for knowledge intensive language tasks, 2021.
- 998  
999 Ronak Pradeep, Sahel Sharifymoghaddam, and Jimmy Lin. Rankvicuna: Zero-shot listwise document  
1000 reranking with open-source large language models, 2023a.
- 1001  
1002 Ronak Pradeep, Sahel Sharifymoghaddam, and Jimmy Lin. Rankzephyr: Effective and robust  
1003 zero-shot listwise reranking is a breeze!, 2023b.
- 1004  
1005 Yifu Qiu, Hongyu Li, Yingqi Qu, Ying Chen, Qiaoqiao She, Jing Liu, Hua Wu, and Haifeng Wang.  
1006 Dureader\_retrieval: A large-scale chinese benchmark for passage retrieval from web search engine,  
1007 2022.
- 1008  
1009 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language  
1010 understanding by generative pre-training, 2018. URL [https://cdn.openai.com/  
1011 research-covers/language-unsupervised/language\\_understanding\\_  
1012 paper.pdf](https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf).
- 1013  
1014 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever,  
1015 et al. Language models are unsupervised multitask learners, 2019. URL [https:  
1016 //d4mucfpksywv.cloudfront.net/better-language-models/language\\_  
1017 models\\_are\\_unsupervised\\_multitask\\_learners.pdf](https://d4mucfpksywv.cloudfront.net/better-language-models/language_models_are_unsupervised_multitask_learners.pdf).
- 1018  
1019 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
1020 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.  
1021 Learning transferable visual models from natural language supervision, 2021.
- 1022  
1023 Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John  
1024 Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. Scaling language models:  
1025 Methods, analysis & insights from training gopher, 2022.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2023.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer, 2023.

- 1026 Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for  
1027 machine comprehension of text, 2016.  
1028
- 1029 Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel  
1030 Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom Eccles, Jake  
1031 Bruce, Ali Razavi, Ashley Edwards, Nicolas Heess, Yutian Chen, Raia Hadsell, Oriol Vinyals,  
1032 Mahyar Bordbar, and Nando de Freitas. A generalist agent, 2022.
- 1033 Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks,  
1034 2019.  
1035
- 1036 Ohad Rubin, Jonathan Herzig, and Jonathan Berant. Learning to retrieve prompts for in-context  
1037 learning, 2022.
- 1038 Alexander M. Rush, Sumit Chopra, and Jason Weston. A neural attention model for abstractive  
1039 sentence summarization, 2015.  
1040
- 1041 Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai,  
1042 Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. Multitask prompted training  
1043 enables zero-shot task generalization, 2022.
- 1044 Teven Le Scao, Thomas Wang, Daniel Hesslow, Lucile Saulnier, Stas Bekman, M Saiful Bari, Stella  
1045 Biderman, Hady Elsahar, Niklas Muennighoff, Jason Phang, Ofir Press, Colin Raffel, Victor Sanh,  
1046 Sheng Shen, Lintang Sutawika, Jaesung Tae, Zheng Xin Yong, Julien Launay, and Iz Beltagy.  
1047 What language model to train if you have one million gpu hours?, 2022.
- 1048 Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer,  
1049 Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to  
1050 use tools, 2023.  
1051
- 1052 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
1053 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 1054 Zejiang Shen, Kyle Lo, Lauren Yu, Nathan Dahlberg, Margo Schlanger, and Doug Downey. Multi-  
1055 lexsum: Real-world summaries of civil rights lawsuits at multiple granularities, 2022.  
1056
- 1057 Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettle-  
1058 moyer, and Wen tau Yih. Replug: Retrieval-augmented black-box language models, 2023.
- 1059 Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus  
1060 Rohrbach, and Douwe Kiela. Flava: A foundational language and vision alignment model, 2022.  
1061
- 1062 Shivalika Singh, Freddie Vargus, Daniel Dsouza, Börje F. Karlsson, Abinaya Mahendiran, Wei-  
1063 Yin Ko, Herumb Shandilya, Jay Patel, Deividas Mataciunas, Laura OMahony, Mike Zhang,  
1064 Ramith Hettiarachchi, Joseph Wilson, Marina Machado, Luisa Souza Moura, Dominik Krzemiński,  
1065 Hakimeh Fadaei, Irem Ergün, Ifeoma Okoh, Aisha Alaagib, Oshan Mudannayake, Zaid Alyafeai,  
1066 Vu Minh Chien, Sebastian Ruder, Surya Guthikonda, Emad A. Alghamdi, Sebastian Gehrmann,  
1067 Niklas Muennighoff, Max Bartolo, Julia Kreutzer, Ahmet Üstün, Marzieh Fadaee, and Sara Hooker.  
1068 Aya dataset: An open-access collection for multilingual instruction tuning, 2024.
- 1069 Jacob Mitchell Springer, Suhas Kotha, Daniel Fried, Graham Neubig, and Aditi Raghunathan.  
1070 Repetition improves language model embeddings. *arXiv preprint arXiv:2402.15449*, 2024.
- 1071 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam  
1072 Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the  
1073 imitation game: Quantifying and extrapolating the capabilities of language models, 2023.  
1074
- 1075 Hongjin Su, Weijia Shi, Jungo Kasai, Yizhong Wang, Yushi Hu, Mari Ostendorf, Wen tau Yih,  
1076 Noah A. Smith, Luke Zettlemoyer, and Tao Yu. One embedder, any task: Instruction-finetuned  
1077 text embeddings, 2023.
- 1078 Ming-Hsiang Su, Chung-Hsien Wu, Kun-Yi Huang, Qian-Bei Hong, and Hsin-Min Wang. A chatbot  
1079 using lstm-based multi-layer embedding for elderly care, 2017. URL [https://ieeexplore.  
ieee.org/document/8336091](https://ieeexplore.ieee.org/document/8336091).

- 1080 Weiwei Sun, Lingyong Yan, Xinyu Ma, Shuaiqiang Wang, Pengjie Ren, Zhumin Chen, Dawei Yin,  
1081 and Zhaochun Ren. Is chatgpt good at search? investigating large language models as re-ranking  
1082 agents, 2023.
- 1083 Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks,  
1084 2014.
- 1085 Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung,  
1086 Aakanksha Chowdhery, Quoc V. Le, Ed H. Chi, Denny Zhou, and Jason Wei. Challenging  
1087 big-bench tasks and whether chain-of-thought can solve them, 2022.
- 1088 Flax Sentence Embeddings Team. Stack exchange question pairs, 2021a. URL [https://hf.co/  
1089 datasets/flax-sentence-embeddings/](https://hf.co/datasets/flax-sentence-embeddings/).
- 1090 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu  
1091 Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, et al. Gemini: A family of highly capable  
1092 multimodal models, 2023.
- 1093 Sentence Transformers Team. (title, body) pairs from the npr.org website, 2021b. URL [https:  
1094 //hf.co/datasets/sentence-transformers/embedding-training-data](https://hf.co/datasets/sentence-transformers/embedding-training-data).
- 1095 Sentence Transformers Team. Reddit title body, 2021c. URL [https://hf.co/datasets/  
1096 sentence-transformers/reddit-title-body](https://hf.co/datasets/sentence-transformers/reddit-title-body).
- 1097 Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. Beir: A  
1098 heterogenous benchmark for zero-shot evaluation of information retrieval models, 2021.
- 1099 James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. Fever: a large-scale  
1100 dataset for fact extraction and verification, 2018.
- 1101 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
1102 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand  
1103 Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language  
1104 models, 2023.
- 1105 Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada,  
1106 Shengyi Huang, Leandro von Werra, Clémentine Fourier, Nathan Habib, Nathan Sarrazin, Omar  
1107 Sanseviero, Alexander M. Rush, and Thomas Wolf. Zephyr: Direct distillation of lm alignment,  
1108 2023.
- 1109 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz  
1110 Kaiser, and Illia Polosukhin. Attention is all you need, 2023.
- 1111 Ben Wang and Aran Komatsuzaki. Gpt-j-6b: A 6 billion parameter autoregressive language model,  
1112 2021. URL <https://github.com/kingoflolz/mesh-transformer-jax>.
- 1113 Bin Wang and C. C. Jay Kuo. Sbert-wk: A sentence embedding method by dissecting bert-based  
1114 word models, 2020.
- 1115 Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder,  
1116 and Furu Wei. Text embeddings by weakly-supervised contrastive pre-training, 2022a.
- 1117 Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. Improving  
1118 text embeddings with large language models, 2024.
- 1119 Thomas Wang, Adam Roberts, Daniel Hesslow, Teven Le Scao, Hyung Won Chung, Iz Beltagy,  
1120 Julien Launay, and Colin Raffel. What language model architecture and pretraining objective work  
1121 best for zero-shot generalization?, 2022b.
- 1122 Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei,  
1123 Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, et al.  
1124 Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks, 2022c.

- 1134 Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Raghavi Chandu,  
1135 David Wadden, Kelsey MacMillan, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. How far  
1136 can camels go? exploring the state of instruction tuning on open resources, 2023.  
1137
- 1138 Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du,  
1139 Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners, 2022.  
1140
- 1141 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le,  
1142 and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023.  
1143
- 1144 Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. Sheared llama: Accelerating language  
1145 model pre-training via structured pruning, 2023.  
1146
- 1147 Shitao Xiao and Zheng Liu. Retromae v2: Duplex masked auto-encoder for pre-training retrieval-  
1148 oriented language models, 2022.  
1149
- 1149 Shitao Xiao, Zheng Liu, Yingxia Shao, and Zhao Cao. Retromae: Pre-training retrieval-oriented  
1150 language models via masked auto-encoder, 2022.  
1151
- 1151 Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. C-pack: Packaged resources to  
1152 advance general chinese embedding, 2023.  
1153
- 1154 Xiaohui Xie, Qian Dong, Bingning Wang, Feiyang Lv, Ting Yao, Weinan Gan, Zhijing Wu, Xiang-  
1155 sheng Li, Haitao Li, Yiqun Liu, and Jin Ma. T2ranking: A large-scale chinese benchmark for  
1156 passage ranking, 2023.  
1157
- 1157 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov,  
1158 and Christopher D. Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question  
1159 answering, 2018.  
1160
- 1161 Michihiro Yasunaga, Armen Aghajanyan, Weijia Shi, Rich James, Jure Leskovec, Percy Liang, Mike  
1162 Lewis, Luke Zettlemoyer, and Wen tau Yih. Retrieval-augmented multimodal language modeling,  
1163 2023.  
1164
- 1164 Jun Yin, Xin Jiang, Zhengdong Lu, Lifeng Shang, Hang Li, and Xiaoming Li. Neural generative  
1165 question answering, 2016.  
1166
- 1167 Zheng-Xin Yong, Hailey Schoelkopf, Niklas Muennighoff, Alham Fikri Aji, David Ifeoluwa Adelani,  
1168 Khalid Almubarak, M Saiful Bari, Lintang Sutawika, Jungo Kasai, Ahmed Baruwa, Genta Indra  
1169 Winata, Stella Biderman, Edward Raff, Dragomir Radev, and Vassilina Nikoulina. Bloom+1:  
1170 Adding language support to bloom for zero-shot prompting, 2023.  
1171
- 1172 Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual  
1173 denotations: New similarity metrics for semantic inference over event descriptions, 2014. URL  
1174 <https://aclanthology.org/Q14-1006>.
- 1175 Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu.  
1176 Coca: Contrastive captioners are image-text foundation models, 2022.  
1177
- 1178 Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text  
1179 classification, 2016.  
1180
- 1180 Xin Zhang, Zehan Li, Yanzhao Zhang, Dingkun Long, Pengjun Xie, Meishan Zhang, and Min Zhang.  
1181 Language models are universal embedders, 2023.  
1182
- 1183 Xinyu Zhang, Xueguang Ma, Peng Shi, and Jimmy Lin. Mr. tydi: A multi-lingual benchmark for  
1184 dense retrieval, 2021.  
1185
- 1186 Xinyu Zhang, Nandan Thakur, Odunayo Ogundepo, Ehsan Kamalloo, David Alfonso-Hermelo,  
1187 Xiaoguang Li, Qun Liu, Mehdi Rezagholizadeh, and Jimmy Lin. Making a mirac: Multilingual  
information retrieval across a continuum of languages, 2022.

1188 Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright, Hamid  
1189 Shojanazeri, Myle Ott, Sam Shleifer, Alban Desmaison, Can Balioglu, Pritam Damania, Bernard  
1190 Nguyen, Geeta Chauhan, Yuchen Hao, Ajit Mathews, and Shen Li. Pytorch fsdp: Experiences on  
1191 scaling fully sharded data parallel, 2023.

1192 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
1193 Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.  
1194 Judging llm-as-a-judge with mt-bench and chatbot arena, 2023.

1195  
1196 Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat,  
1197 Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy.  
1198 Lima: Less is more for alignment, 2023.

1199  
1200 Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng, Haonan Chen,  
1201 Zhicheng Dou, and Ji-Rong Wen. Large language models for information retrieval: A survey,  
1202 2024.

1203 Terry Yue Zhuo, Armel Zebaze, Nitchakarn Suppattarachai, Leandro von Werra, Harm de Vries, Qian  
1204 Liu, and Niklas Muennighoff. Astraios: Parameter-efficient instruction tuning code large language  
1205 models, 2024.

1206 Ahmet Üstün, Viraat Aryabumi, Zheng-Xin Yong, Wei-Yin Ko, Daniel D’souza, Gbemileke Onilude,  
1207 Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, Freddie Vargus, Phil Blunsom, Shayne  
1208 Longpre, Niklas Muennighoff, Marzieh Fadaee, Julia Kreutzer, and Sara Hooker. Aya model: An  
1209 instruction finetuned open-access multilingual language model, 2024.

1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

# APPENDIX

## Contents

1242		
1243		
1244		
1245		
1246		
1247		
1248		
1249	<b>A Ablations</b>	<b>25</b>
1250		
1251	<b>B Discussion</b>	<b>28</b>
1252		
1253	<b>C Aligning GRITLM</b>	<b>29</b>
1254		
1255	<b>D Few-shot embedding does not work</b>	<b>30</b>
1256		
1257	<b>E RAG Caching CPU Latency</b>	<b>30</b>
1258		
1259	<b>F Additional RAG results</b>	<b>31</b>
1260		
1261	<b>G Loss Curves</b>	<b>31</b>
1262		
1263	<b>H Evaluation</b>	<b>32</b>
1264		
1265	<b>I Ablations Detailed Results</b>	<b>33</b>
1266		
1267	<b>J GRITLM MTEB Full Results</b>	<b>38</b>
1268		
1269	<b>K Reducing Embedding Training Memory</b>	<b>39</b>
1270		
1271	<b>L Hyperparameters</b>	<b>40</b>
1272		
1273	<b>M Embedding Instruction for Generative Models</b>	<b>40</b>
1274		
1275	<b>N HumanEval Format</b>	<b>41</b>
1276		
1277	<b>O Embedding in FP32 vs BF16</b>	<b>41</b>
1278		
1279	<b>P Unreliability of MT-Bench</b>	<b>42</b>
1280		
1281	<b>Q Limitations and Future Work</b>	<b>42</b>
1282		
1283	<b>R Dataset Composition</b>	<b>43</b>
1284		
1285	<b>S Dataset Samples</b>	<b>45</b>
1286		
1287	<b>T Evaluation Prompts</b>	<b>48</b>
1288	T.1 Embedding Prompts . . . . .	48
1289	T.2 Embedding Few-Shot Prompts . . . . .	54
1290	T.3 Generative Prompts . . . . .	58
1291	T.4 RAG Prompts . . . . .	64
1292		
1293	<b>U Hardware</b>	<b>67</b>
1294		
1295	<b>V Artifacts</b>	<b>67</b>



## 1296 A ABLATIONS

1297  
1298 **Attention and pooling** We train GRITLM starting from a pretrained decoder language model  
1299 which has been trained with causal attention. Prior work has shown that while embeddings of causal  
1300 LLMs are competitive, they are outperformed by BERT-like encoders with bidirectional attention at  
1301 the same number of parameters (Muennighoff, 2022; Devlin et al., 2019). This lines up with intuition,  
1302 as bidirectional attention allows the model to adjust the representation of the first tokens based on  
1303 information obtained from future tokens. Meanwhile, causal attention only allows information to  
1304 propagate one way. Thus, for causal attention early tokens may yield poor representations due to a lack  
1305 of understanding of the entire sample. To counter this issue, we experiment with adapting the model  
1306 during finetuning to learn to use bidirectional attention. In Table 5 we find that **adapting the causally**  
1307 **pretrained LLM with bidirectional attention provides the best embedding performance**. For  
1308 fully causal embeddings, we confirm findings from Muennighoff (2022) that position-weighted mean  
1309 pooling (“Wmean”) leads to better embedding performance than taking the embedding of the last  
1310 token despite recent work finding the opposite (Zhang et al., 2023; Ma et al., 2023a). For last token  
1311 pooling, we follow Zhang et al. (2023) and use a special token. We find that adapting the model to be  
1312 a PrefixLM (Raffel et al., 2023), whereby the attention over the generative instruction is bidirectional  
1313 but still causal for the response (“Sample”) worsens performance in contrast to prior work (Wang et al.,  
1314 2022b). Thus, we stick with fully causal generation. The unified variant significantly outperforms the  
1315 embedding-only variants, while underperforming the best generative-only variant. However, once  
1316 we switched from MEDI to the E5 dataset in later ablations the embedding-only variant matched  
1317 the unified variant. Meanwhile, the worse generative performance of the unified model was due to a  
1318 suboptimal loss setting that we fixed in the loss ablations. Several papers after the initial preprint  
1319 release of this work have confirmed the benefit of bidirectional attention (BehnamGhader et al., 2024;  
1320 Springer et al., 2024).

1321 **Base model** The GRITLM approach generalizes to any generative language model, thus we ablate  
1322 initializing from GPT-J 6B (Wang & Komatsuzaki, 2021), Llama 2 7B or Mistral 7B (Jiang et al.,  
1323 2023). Using Mistral 7B leads to the best performance for both embedding and generative tasks.  
1324 For generative tasks, this is expected as the pretrained Mistral 7B performs the best among the  
1325 three (Table 2). However, for embedding tasks, GPT-J outperforms Mistral 7B (Table 1). Thus, **the**  
1326 **embedding performance of a pretrained model is not predictive of its embedding performance**  
1327 **after finetuning**. Rather, its generative performance appears to be a more reliable indicator of its  
1328 embedding performance after finetuning.

1329 **Generative dataset** We benchmark our filtered Tulu 2 introduced in §3.1 (Iverson et al., 2023)  
1330 with UltraChat (Ding et al., 2023; Tunstall et al., 2023) and the OpenAssistant version from Oc-  
1331 toPack (Muennighoff et al., 2023a; Köpf et al., 2023; Longpre et al., 2023). Using Tulu 2 leads to  
1332 better performance on every generative task considered (see Appendix I for per-task results). This  
1333 is likely due to Tulu 2 containing a larger diversity of tasks (Iverson et al., 2023). Another possible  
1334 reason is that Tulu 2 may have been carefully tuned on the generative evaluation datasets, as we use  
1335 largely the same evaluation setup as the creators of Tulu 2 (Iverson et al., 2023).

1336 **Embedding dataset** We benchmark MEDI (Su et al., 2023), a new version of MEDI with better  
1337 negatives which we build and call MEDI2, and the E5 dataset (Wang et al., 2024). While MEDI  
1338 and MEDI2 always preface instructions with “Represent” (see e.g. Figure 11), the E5 dataset places  
1339 no constraint on the instruction prefix (see e.g. Figure 12). Thus, when using the E5 dataset the  
1340 “<|embed|>” formatting is critical to tell the model that it will be subject to the representation loss,  
1341 not the generative loss (Figure 3). Further, MEDI and MEDI2 always contain instructions for both  
1342 queries and documents, which we refer to as **two-sided instructions**. Meanwhile, the E5 dataset  
1343 uses **one-sided instructions** for asymmetric datasets (Muennighoff, 2022), whereby the documents  
1344 receive no instructions, only the queries. The advantage of not using document instructions is that  
1345 the document corpus can be encoded once and then cached and reused across a variety of tasks.  
1346 During training on E5, symmetric tasks are also in a one-sided setting, but we still evaluate them in  
1347 the two-sided format. This should not be a problem as the cosine similarity function we use during  
1348 training is transitive: if sentence A with instruction is similar to sentence B without instruction, and  
1349 sentence B without instruction is similar to sentence C with instruction, then we can confidently say  
that sentence A with instruction is also similar to sentence C with instruction. As depicted in Table 5,

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

Attention Emb		Attention Gen		Pooling	Emb	Gen
Instruction	Sample	Instruction	Sample			
<i>Embedding Only</i>						
	Causal			Wmean	60.0	-
	Causal Bidirectional			Mean	61.0	-
	Bidirectional			Mean	61.8	-
<i>Generative Only</i>						
		Causal			-	<b>55.2</b>
		Bidirectional	Causal		-	50.7
<i>Unified</i>						
	Causal	Causal		Last token	61.2	53.0
	Causal	Causal		Wmean	62.8	52.8
	<b>Bidirectional</b>	<b>Causal</b>		<b>Mean</b>	<b>64.0</b>	52.9

(a) Attention and pooling ablations. Wmean is position-weighted mean pooling (Muennighoff, 2022).

Variant	Emb	Gen	Dataset	Emb	Dataset	Gen
<b>Mistral 7B</b>	<b>54.6</b>	<b>22.4</b>	MEDI	64.0	<b>Tulu 2</b>	<b>55.2</b>
Llama 2 7B	48.2	20.8	MEDI2	64.7	OASST	37.7
GPT-J 6B	51.9	14.0	<b>E5</b>	<b>66.0</b>	UltraChat	47.4

(b) Base model

(c) Embedding dataset

(d) Generative dataset

Variant	Emb	Gen	BS Emb:Gen	Emb	Gen	Precision	Emb	Gen
<b>No head</b>	<b>62.7</b>	<b>49.2</b>	256:256	63.2	<b>53.4</b>	FP32	66.3	52.4
→ 1024	62.1	48.0	4096:256	<b>64.2</b>	53.3	<b>BF16*</b>	<b>66.5</b>	<b>55.0</b>

(e) Embedding head

(f) Batch size (BS)

(g) Precision

IBN origin	Emb	Gen	Format	Gen	Tokens	Emb	Gen
Any dataset	<b>66.0</b>	50.9	<b>Tulu 2</b>	<b>55.2</b>	512	64.1	52.2
<b>Same dataset</b>	<b>66.0</b>	<b>51.1</b>	Zephyr $\beta$	49.0	<b>2048</b>	<b>64.7</b>	<b>53.8</b>

(h) In-batch negatives (IBN)

(i) Format

(j) Emb training max tokens

Gen loss type	$\mathcal{L}_{Rep}/\mathcal{L}_{Gen}$	Emb	Gen	Gen loss type	AlpacaEval
Token	2.4	66.1	54.4	Mix (4 → 64)	67.6
Token	6.0	66.5	55.0	<b>Mix (32 → 8)</b>	<b>74.7</b>
<b>Mix (32 → 8)</b>	4.1	<b>66.7</b>	<b>55.4</b>		

(k) Loss ablations.  $\mathcal{L}_{Rep}/\mathcal{L}_{Gen}$  is the loss ratio of the 1st step adjusted via  $\lambda_{Rep}$  and  $\lambda_{Gen}$ . Mix refers to mixing sample and token level loss, e.g. (32 → 8) is token level loss across 32 samples and then sample level loss across 8 sub-batches for a total batch size of 256.

Table 5: **GRIT ablations**. Emb corresponds to the MTEB average, while Gen corresponds to the average across generative tasks (Appendix H). The embedding head variant “→ 1024” corresponds to down-projecting the final hidden state with a linear layer from 4096 to 1024 dimensions, only for embedding tasks. BF16\* means that some computations are still in FP32 as explained in Appendix A. The setting chosen for GRITLM is **bold**. Once an ablation was successful, we adopted its setting, thus the bold performance slightly varies from one table to the next. For example, the base model ablation (b) is done for just 100 hundred steps with sub-optimal formatting. Full results are in Appendix I.

1404 using the E5 dataset performs best by a wide margin. An inspection of samples, suggests that this is  
1405 likely due to its superior hard negatives and diversity of tasks generated by GPT-4 (Appendix S). For  
1406 our final runs with the E5 dataset, we additionally add scientific data (§3.1).  
1407

1408 **Embedding head** The cost of caching the embeddings of a large document corpus is directly  
1409 proportional to the embedding dimensionality. To minimize such costs, we experiment with adding  
1410 an embedding head consisting of a linear layer with activation that down-projects the embedding (Ni  
1411 et al., 2021a; Muennighoff, 2022). This layer is only used for embedding tasks. Down-projecting  
1412 the embeddings four-fold (from 4096 to 1024) leads to an embedding performance decrease of  
1413 around 1%. This may be acceptable for certain use cases where the saved storage is more important.  
1414 However, for our final model, we do not use such a head to keep it simple and achieve maximum  
1415 performance. Search techniques (Arya et al., 1998; Johnson et al., 2017; Douze et al., 2024) or  
1416 dimensionality reduction techniques such as Principal Component Analysis still allow for reducing  
1417 the embedding dimension of our final model post-training while maintaining most of the performance.  
1418 Similar to the storage cost-performance trade-off we explore here, we hypothesize that there is a  
1419 speed/cost-performance trade-off with taking the embedding from different layers of our model. For  
1420 example, we could train using the embedding after half the layers of the model, thus speeding up the  
1421 embedding model by 50% while likely only incurring a small drop in embedding performance

1422 **Batch size** Due to the utilization of in-batch negatives for contrastive training (§2), a larger batch  
1423 size provides a more accurate gradient. Thus, scaling up the batch size is a key ingredient in most  
1424 well-performing embedding models (Xiao et al., 2023; Wang et al., 2022a). We experiment with  
1425 scaling up the embedding batch size to 4096 while keeping it at 256 for generative data. This leads to  
1426 a 1.0 gain on the embedding average while generative performance remains stable. Especially the 15  
1427 retrieval datasets that are part of the embedding average benefit from the increase in batch size (see  
1428 Table 18). For our final model, we use a batch size of 2048 for embedding and 256 for generative  
1429 data.

1430 **Precision** The parameters of the Mistral 7B model are in bfloat16 (BF16) precision as it was  
1431 pretrained in this format. We experiment with finetuning it with float32 (FP32) precision versus  
1432 keeping the BF16 format and training with mixed precision. FP32 training is more costly, however,  
1433 the additional precision may result in a better model. Our intuition is that more precision is important  
1434 for embedding but not as much for generation. This is because while for generative tasks evaluated  
1435 greedily, the model output is a discretionary argmax over the predictions of the language modeling  
1436 head, for embedding tasks it is a continuous representation. Thus, small differences due to a lack  
1437 of precision may not change the model’s generation but will affect its representation. Hence, for  
1438 embedding tasks, we always cast the hidden states to FP32 during the pooling operation and keep  
1439 them this way for the similarity computation. Not keeping them in FP32 after pooling worsens  
1440 performance slightly, but may be necessary for cheap storage (see Appendix O). In addition, some  
1441 operations such as layer normalization (Ba et al., 2016) are also performed in FP32 even for BF16  
1442 training due to PyTorch autocast (Zhao et al., 2023). In Table 5, we find that there is no benefit from  
1443 doing even more computations in FP32 besides the ones listed above. Thus, we train and evaluate all  
1444 our other models in BF16 mixed precision to speed up training and inference.

1445 **In-batch negatives** We always use in-batch negatives for embedding training (§2), however, we  
1446 ablate whether or not they come from the same dataset. We hypothesize that making them all come  
1447 from the same dataset leads to better negatives as the model needs to distinguish them based on more  
1448 nuanced differences. In practice, we find that the average embedding performance remains around  
1449 the same. However, we notice a 1.3 jump on the 15-dataset Retrieval average (Table 20). Thus, we  
1450 stick with the variant where in-batch negatives stem from the same dataset.  
1451

1452 **Format** Our chosen format is depicted in Figure 3, which is equivalent to Tulu 2 (Iverson et al., 2023)  
1453 for generative tasks. We also benchmark the Zephyr  $\beta$  format (Tunstall et al., 2023), which has an  
1454 additional end-of-sequence token (“</s>”) after each user utterance. We find that it performs worse  
1455 on generative tasks. The additional end-of-sequence after the user utterance increases the likelihood  
1456 of the model generating another end-of-sequence token earlier than necessary. This significantly  
1457 harms HumanEvalSynthesize performance and slightly reduces AlpacaEval, where long generations  
can be critical (see Appendix I for task-specific performance).

**Max tokens** Our base model, Mistral 7B, can handle sequences of arbitrary length due to its sliding window attention (Jiang et al., 2023). As finetuning with longer sequences is more expensive we ablate its benefits. We compare training with a maximum token limit of 512 versus 2048 for embedding documents. For embedding queries, we always use 256, and for generative data, we always use 2048. We find that increasing the embedding document sequence length during training slightly boosts performance on both embedding and generation even though we still evaluate embedding tasks with 512. This boost likely comes from our training data containing many documents beyond 512 tokens, which need to be truncated if the maximum sequence length is 512. Such truncation may remove the critical parts that make two texts a positive or a negative contrastive pair and thus hinder learning. As our embedding evaluation (MTEB) contains few documents longer than 512 tokens there is little truncation happening at evaluation (Muennighoff et al., 2023c; Günther et al., 2024; 2023). Note that just like their base models, our final models GRITLM 7B and GRITLM 8x7B can produce embeddings for sequences of arbitrary length. However, due to a lack of benchmarks, we do not know how well the embeddings of our models perform for input sequences longer than 512 tokens.

**Loss ablations** As detailed in §2, we experiment with both token and sample level generative loss. Further, we ablate the representation and generative loss weights,  $\lambda_{\text{Rep}}$  and  $\lambda_{\text{Gen}}$ . For the unified visual model CoCa, the authors find that giving a weight of 2 to generation and 1 to embedding boosts performance on both streams (Yu et al., 2022). However, rather than the weights, we argue that the loss ratio,  $\mathcal{L}_{\text{Rep}}/\mathcal{L}_{\text{Gen}}$ , is of more interest as it reveals which objective has a larger impact on the optimization of the model. We maintain a ratio of  $\mathcal{L}_{\text{Rep}}/\mathcal{L}_{\text{Gen}} \approx 1$  i.e. giving more weight to the representation loss. This is because the model has already been pretrained with the generative loss, thus we expect less additional generative training to be necessary. Meanwhile, the contrastive loss for embedding data is new to the model, thus we expect more learning to be needed on the embedding side. Further, the embedding loss drops off extremely quickly as can be seen in the loss graphs in Appendix G. Thus, even though the representation loss has a higher weight at the start, throughout training they have very similar weights with both hovering around a loss of 1.0. We find that mixing sample and token level generative loss leads to the best performance by a small margin. As expected in §2, token level loss to some degree is critical for good performance on AlpacaEval. For “Mix (4 -> 64)” token level loss is applied across only 4 samples and then sample level loss across 64 sub-batches, which leads to a 7-point drop in AlpacaEval performance. This drop is accompanied by a decrease in median AlpacaEval generation length from 941 to 865. Thus, token level loss across many samples is critical to maintaining long generations, which directly impacts the AlpacaEval score.

## B DISCUSSION

**Further unification** To the best of our knowledge, GRITLM is the first model to unify text embedding and generation, and thus all text-based language problems, into a single model at strong performance. However, many adjacent directions remain to be improved or unified. **(a) Multilinguality:** Our model is also capable of embedding and generation in non-English languages as seen in its TyDi QA performance (Table 2). However, major performance gains on non-English tasks are likely possible through both data (Muennighoff et al., 2023d; Yong et al., 2023) and architecture changes (Chen et al., 2024a; Feng et al., 2022; Duquenne et al., 2023) targeting multilinguality. **(b) Multimodality:** Many embedding and generative problems are not purely text-based, such as joint embedding of images and text (Radford et al., 2021), generative image captioning (Hossain et al., 2018), image-text pair classification (Muennighoff, 2020; Kiela et al., 2021) or speech versions of every text problem (Kamath et al., 2019). It remains to be explored whether they can be as easily unified as text embedding and generation in this work.

**Why does GRIT work?** GRIT unifies embedding and generative tasks into a single model at no performance loss on either one, which may seem surprising. When the embedding dataset is MEDI2, we show that embedding performance even improves once the generative objective is added compared to an otherwise equivalent embedding-only model (Appendix A). We think that our results confirm that generative language modeling and text embeddings are two sides of the same coin. Both tasks require a model to have a deep understanding of natural language and only differ in the way that understanding is expressed. Possibly, our unified model contains a small number of parameters that act as a switch to make the final representations either useful for mean pooling and subsequent

embedding tasks or primed for the language modeling head and subsequent generative tasks. We are excited about future work exploring what is happening inside of GRITLM. To support such research, we release all our work freely.

**Optimizing RAG with GRITLM** RAG and the caching variants we have presented in this work operate on a frozen language model. Meanwhile, there has been extensive work on optimizing a generative model specifically for interaction with a retrieval system (Gao et al., 2024; Zhu et al., 2024; Asai et al., 2023a). These works commonly optimize only the retriever (Shi et al., 2023) or only the reader (Borgeaud et al., 2022; Yasunaga et al., 2023; Asai et al., 2023b; Luo et al., 2023). However, recent work has shown that jointly optimizing both models leads to the best performance (Lin et al., 2023). With its state-of-the-art retrieval and generative performance, GRITLM can act as both the retriever and reader in a single model. Thus, optimizing either one also changes the parameters of the other. This has the potential to significantly simplify the joint optimization of the retriever and reader. For example, it may suffice to only use the next-token objective (Equation 2) to penalize the retriever for providing irrelevant context and at the same time the reader for poor use of the given context. This is in contrast to separate models and objective functions used in Lin et al. (2023).

## C ALIGNING GRITLM

Table 6: **Aligning GRITLM with KTO after GRIT.** The upper table depicts embedding performance while the lower depicts generative performance.

Task (→)	CLF	Clust.	PairCLF	Rerank	Retrieval	STS	Summ.	Avg.
Metric (→)	Acc.	V-Meas.	AP	MAP	nDCG	Spear.	Spear.	
Dataset # (→)	12	11	3	4	15	10	1	56
GRITLM 7B	79.5	50.6	87.2	60.5	57.4	83.4	30.4	66.8
GRITLM 7B KTO	79.6	50.1	87.1	60.5	57.1	83.5	30.5	66.7
GRITLM 8x7B	78.5	50.1	85.0	59.8	55.1	83.3	29.8	65.7
GRITLM 8x7B KTO	78.7	50.0	84.4	59.4	54.1	82.5	30.8	65.2
Dataset (→)	MMLU	GSM8K	BBH	TyDi QA	HumanEval	Alpaca	Avg.	
Setup (→)	0 FS	8 FS, CoT	3 FS, CoT	1 FS, GP	0 FS	0 FS, 1.0		
Metric (→)	EM	EM	EM	F1	pass@1	% Win		
GRITLM 7B	57.6	57.5	54.8	55.4	32.8	74.8	55.5	
GRITLM 7B KTO	57.6	57.5	55.4	55.8	31.5	86.7	57.4	
GRITLM 8x7B	66.7	61.5	70.2	58.2	53.4	84.0	65.7	
GRITLM 8x7B KTO	66.8	79.5	67.1	31.4	56.8	95.3	66.2	

It is common to follow the instruction finetuning stage of generative language models by an alignment tuning stage using methods like PPO (Schulman et al., 2017), DPO (Rafailov et al., 2023), or KTO (Ethayarajh et al., 2024) (“HALOs” (Ethayarajh et al., 2024)). We experiment with further finetuning GRITLM using KTO and benchmark the resulting models in Table 6. During this KTO stage, no further embedding training is performed, thus it leads to a slight performance drop on the MTEB average (66.8 to 66.7 and 65.7 to 65.2). However, the average generative performance of the KTO-tuned models is stronger. Notably, AlpacaEval jumps by  $\approx 10$  points for both models. On the more recent Alpaca 2.0 (Dubois et al., 2024), GritLM-8x7B-KTO has a length-controlled win rate of 18.5 with an average length of 1662 (not depicted). Thus, the KTO-finetuned models may be more useful for use cases where the generative performance is more important. Future work may consider continuing the embedding training during the alignment tuning stage. It may also be possible to develop an alignment tuning method specifically for embedding performance and combine it with generative alignment via KTO.

### D FEW-SHOT EMBEDDING DOES NOT WORK

For generative models it has been well-established that providing in-context examples (“few-shots”, FS) improves performance (Brown et al., 2020). However, to the best of our knowledge, there has been no work on in-context learning with embedding models. In Table 7, we benchmark the default 0-shot format with providing a single few-shot example following the task instruction. We take the few-shot example from the respective evaluation dataset (see §T.2 for the prompts). We find that providing few-shot examples overall worsens performance. While there are small gains among PairClassification tasks (SprintDup. and TwitterURL), these are marginal and inconsistent. For the model trained on MEDI2, we even include few-shot embedding samples in the training data for around 5% of training samples. However, the model seems not to have learned to make good use of the few-shot examples.

Table 7: **Few-shot embedding.** The 12 MTEB datasets (“DS”) are grouped by the 7 main MTEB tasks in the same order as in Table 1.

Train DS (→) MTEB DS (↓)	E5S		MEDI2	
	0 FS	1 FS	0 FS	1 FS
Banking77	<b>88.5</b>	88.3	<b>88.1</b>	87.9
Emotion	<b>52.8</b>	51.0	<b>52.5</b>	51.9
IMDB	<b>95.0</b>	93.9	<b>94.3</b>	92.2
BiorxivS2S	<b>39.8</b>	39.4	<b>37.6</b>	37.4
SprintDup.	93.0	<b>94.9</b>	95.2	<b>95.7</b>
TwitterSem	<b>81.1</b>	77.9	<b>76.8</b>	73.9
TwitterURL	<b>87.4</b>	87.1	85.9	<b>86.1</b>
ArguAna	<b>63.2</b>	51.7	<b>53.5</b>	53.2
SCIDOCS	<b>24.4</b>	19.7	<b>25.5</b>	25.5
AskUbuntu	<b>67.3</b>	64.7	<b>66.6</b>	66.0
STS12	77.3	<b>78.0</b>	<b>76.6</b>	73.5
SummEval	<b>30.4</b>	29.5	29.1	<b>31.5</b>

### E RAG CACHING CPU LATENCY

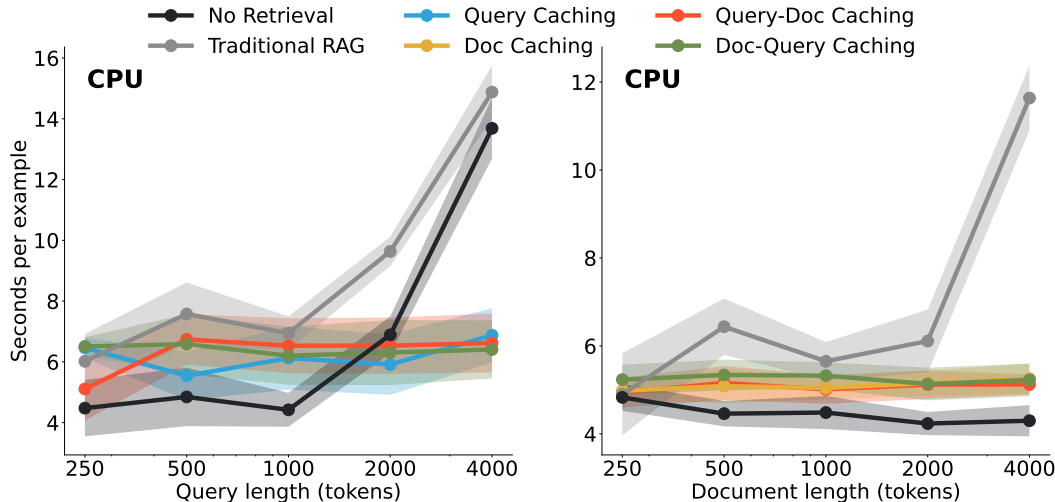


Figure 6: **Inference latency of RAG with GRITLM 7B on CPUs.** When benchmarking scaling query length (left), document length is fixed at 1, whereas query length is fixed at 1 when scaling document length (right). In addition to the query/doc lengths, the formatting and prompt take up around 40 tokens. We visualize the standard deviation across 100 runs as the shaded area. For each approach, we generate 16 tokens. See Figure 5 for GPU latency.

## F ADDITIONAL RAG RESULTS

In §5, we find that doc caching is the most promising caching variant out of the ones we propose. This is because (a) documents are usually significantly longer than queries, thus caching documents has the highest potential to reduce latency, (b) it maintains performance of regular RAG (Table 4, and (c) it even works for non-GRIT models though it requires more time to construct the cache for non-GRIT models (§5). Thus we further experiment with doc caching in Table 8 to verify its performance on other datasets. Similar to Natural Questions in Table 4, we observe that doc caching maintains performance of regular RAG (even slightly improves) for TriviaQA and MMLU despite the attention mismatch. Note that the attention mismatch problem can always be resolved by simply not using bidirectional attention for the embedding part and thereby guarantee the same performance as not using RAG, however, not using bidirectional attention comes at a slight reduction in embedding performance according to our ablation experiments (Appendix A).

We also benchmark the BGE series of embedding models (Xiao et al., 2023) in Table 9 for RAG. We find performance to be significantly worse than with GRITLM in Table 4. Based on a manual inspection of samples, it appears that the embedding models commonly retrieve irrelevant passages that confuse the generative model. There may be other smaller embedding models or other generative models that may perform better, but overall we expect the RAG performance to be a function of the embedding and generative performance of the individual components (e.g. if an embedding model performs better than GRITLM, we would expect it to lead to better RAG performance; BGE generally does not perform better on embedding as shown in Table 1).

Table 8: **Additional doc caching results.** We use the same setup as in Table 4 to benchmark doc caching on two additional datasets: TriviaQA (Joshi et al., 2017) and MMLU (Hendrycks et al., 2022).

Dataset (→) Metric (→)	TriviaQA Match (0-shot, ↑)	MMLU Match (0-shot, ↑)
RAG	52.12	51.10
Doc Caching	<b>57.93</b>	<b>53.46</b>

Table 9: **Additional RAG results with BGE.** We use the same setup as in Table 4 to benchmark BGE embedding models with the “Query then document” prompt. The generative model is still GRITLM 7B.

Dataset (→) Metric (→)	NQ Match (0-shot, ↑)
BGE Large 0.34B	10.39
BGE Base 0.11B	10.31
BGE Small 0.03B	10.17

## G LOSS CURVES

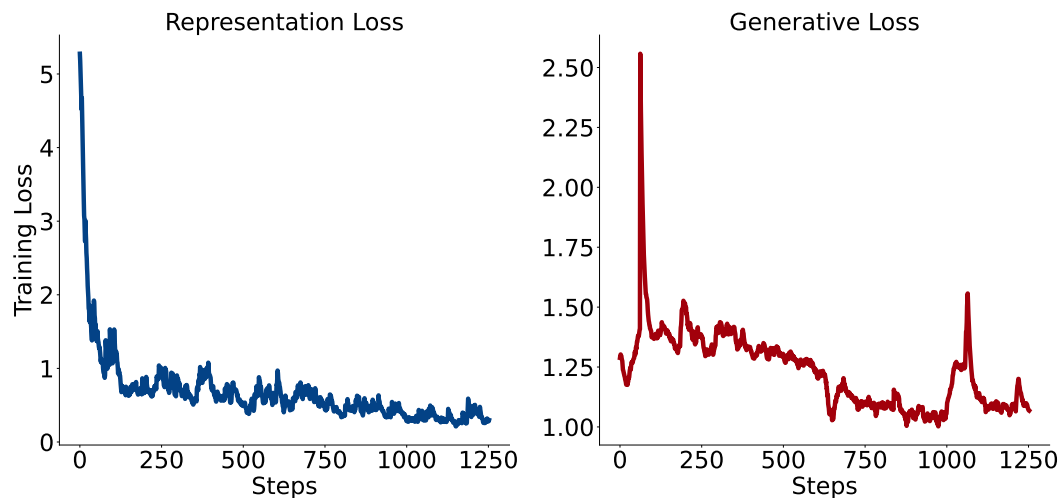


Figure 7: **GRITLM 7B training loss smoothed with exponential moving average smoothing and a weight of 0.9.**

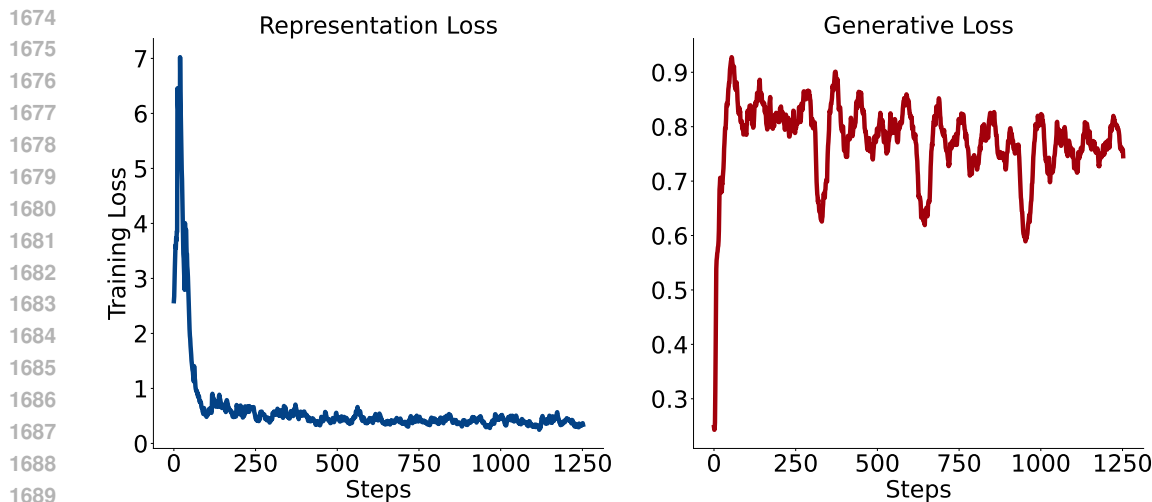


Figure 8: GRITLM 8x7B training loss smoothed with exponential moving average smoothing and a weight of 0.9.

## H EVALUATION

For evaluating GRITLM, we select the most commonly used embedding and generative benchmarks:

**Embedding** To evaluate embedding performance we use the 7 main tasks from MTEB (Muenighoff et al., 2023c).

(1) **Classification (CLF)**: A logistic regression classifier is trained on embeddings from texts with different labels. The classifier is scored with **F1**.

(2) **Clustering (Clust.)**: K-means clustering is performed on embeddings from different sources. The agreement of the clusters with respect to the source labels is scored with **V-measure**.

(3) **Pair Classification (PairCLF)**: The cosine similarity of two embeddings with a binary label is computed. The optimal similarity threshold across all samples is found and scored with **AP** (average precision).

(4) **Reranking (Rerank)** A query embedding and reference embeddings are compared with cosine similarity. The similarities are scored versus the ground truth ranking of the references via **MAP** (mean AP).

(5) **Retrieval**: A query embedding and embeddings of references are compared with cosine similarity. The position of the correct reference(s) in the top ten with the highest cosine similarity is scored with **nDCG@10** (normalized discounted cumulative gain).

(6) **STS**: The cosine similarity of two embeddings is compared with a ground truth continuous score of their similarity and scored with **Spearman** correlation.

(7) **Summarization (Summ.)** Human-written and machine-written summaries of the same text are embedded. The cosine similarity of the embeddings is compared to human ratings of the machine summaries and scored with **Spearman** correlation.

Among the tasks, Reranking, Retrieval, and Summarization are asymmetric i.e. there are two different kinds of embeddings: queries and documents. Others are symmetric i.e. there is only one kind. We use instructions for every dataset specified in §T.1. Notably, for some models, we use different instructions for query and document embeddings when dealing with asymmetric tasks. The datasets within each task cover diverse domains ranging from scientific papers to casual conversations.

**Generation** For evaluating the generative performance of GRITLM, we largely follow the evaluation setup of Tülu (Wang et al., 2023; Ivison et al., 2023) using open-source frameworks (Gao et al., 2021a; Ben Allal et al., 2022).

(1) **Multiple-Choice Question Answering via MMLU (Hendrycks et al., 2022)**: Models are tasked to answer knowledge-intensive questions from different fields, such as humanities, social sciences, and hard sciences. No few-shots are provided and answers are evaluated with **exact match**.



(2) **Problem solving via GSM (Cobbe et al., 2021)**: Models are tasked to solve a math problem requiring multi-step reasoning. 8 few-shot (FS) examples with chain-of-thought reasoning (CoT) (Wei et al., 2023) are provided and **exact match** is measured.

(3) **Multilingual Closed-book Question Answering via TyDi QA (Clark et al., 2020)**: Models are tasked to answer a question in one of six languages. We evaluate in the Gold Passage and no-context setting following Anil et al. (2023).

(4) **Code Generation via HumanEvalSynthesize (Muennighoff et al., 2023a; Chen et al., 2021)**: We use the HumanEvalSynthesize Python dataset (Muennighoff et al., 2023a), which is adapted from HumanEval (Chen et al., 2021) for easy evaluation of instruction-following models. Using the instruction format is different from Ivison et al. (2023) who use HumanEval without an instruction format which is not how the model is used in practice. Following Muennighoff et al. (2023a), we score pass@1 using 20 samples and a temperature of 0.2.

(5) **Boolean Expressions, Causal Judgement, etc. via BBH (Srivastava et al., 2023; Suzgun et al., 2022)** We evaluate a variety of reasoning tasks using BIG-Bench Hard (BBH) (Srivastava et al., 2023; Suzgun et al., 2022). Similar to GSM8K, 3 FS CoT examples are provided and **exact match** is measured.

(6) **Open-ended writing, Summarization, Role-playing, etc. via AlpacaEval (Alpaca) (Li et al., 2023b; Dubois et al., 2023)** We evaluate a variety of open-ended generation tasks via the original 1.0 version of AlpacaEval (Li et al., 2023b; Dubois et al., 2023). GPT-4 (OpenAI et al., 2023) is used to determine the win rate of generations compared to provided GPT-3 (Brown et al., 2020) answers. We differ from Ivison et al. (2023) in that we reduce the maximum token length to 6144 from 8192. We do not use MT-Bench due to its limitations pointed out in Appendix P. To ensure reproducibility, we use greedy evaluation throughout.

## I ABLATIONS DETAILED RESULTS

We display a breakdown of the results from Table 5 in Table 10 to Table 21. For MTEB per-dataset results, we refer to Appendix J, the MTEB leaderboard (<https://huggingface.co/spaces/mteb/leaderboard>) and our released result files (<https://huggingface.co/datasets/ANONYMIZED>).

Table 10: **Unified models attention and pooling ablations.** The sequence of Cs and Bs refers to the attention mechanism for (from left to right): Emb instruction, Emb sample, Gen instruction, Gen sample, where C=Causal, B=Bidirectional, Emb=Embedding and Gen=Generative. WM, LT and M refer to position-weighted mean, last token and mean pooling, respectively.

Task (→)	CLF	Clust.	PairCLF	Rerank	Retrieval	STS	Summ.	Avg.
Metric (→)	Acc.	V-Meas.	AP	MAP	nDCG	Spear.	Spear.	
Dataset # (→)	12	11	3	4	15	10	1	56
CCCC WM	77.9	47.9	81.5	59.0	49.4	80.3	29.4	62.8
CCCC LT	78.8	46.9	84.5	59.6	43.9	78.7	29.3	61.2
BBCC M	79.0	48.6	86.3	59.5	49.9	81.7	30.1	63.8
Dataset (→)	MMLU	GSM8K	BBH	TyDi QA	HumanEval	Alpaca	Avg.	
Setup (→)	0 FS	8 FS, CoT	3 FS, CoT	1 FS, GP	0 FS	0 FS, 1.0		
Metric (→)	EM	EM	EM	F1	pass@1	% Win		
CCCC WM	57.5	45.0	53.1	56.0	32.3	72.9	52.8	
CCCC LT	57.2	45.5	54.7	54.0	31.1	75.7	53.0	
BBCC M	57.0	46.5	54.5	55.0	30.4	73.8	52.9	

Table 11: **Embedding-only models attention and pooling ablations.** The sequence of Cs and Bs refers to the attention mechanism for (from left to right): Emb instruction, Emb sample, where C=Causal, B=Bidirectional and Emb=Embedding. WM and M refer to position-weighted mean and mean pooling, respectively.

Task (→)	CLF	Clust.	PairCLF	Rerank	Retrieval	STS	Summ.	Avg.
Metric (→)	Acc.	V-Meas.	AP	MAP	nDCG	Spear.	Spear.	
Dataset # (→)	12	11	3	4	15	10	1	56
CC WM	77.1	44.0	83.3	57.0	43.2	79.6	29.4	60.0
CB M	76.4	45.5	83.1	56.8	45.7	80.6	30.4	61.0
BB M	77.3	46.0	83.8	58.2	46.8	81.0	32.3	61.8

Table 12: **Generative-only models attention ablations.** The sequence of Cs and Bs refers to the attention mechanism for (from left to right): Gen instruction, Gen sample, where C=Causal and B=Bidirectional. IL=interleaved, whereby the bidirectional attention is interleaved with causal attention in multi-turn samples (bidirectional for instructions, causal for answers). This allows for faster generation in multi-turn settings as the kv-cache of the answer can be reused.

Dataset (→)	MMLU	GSM8K	BBH	TyDi QA	HumanEval	Alpaca	Avg.
Setup (→)	0 FS	8 FS, CoT	3 FS, CoT	1 FS, GP	0 FS	0 FS, 1.0	
Metric (→)	EM	EM	EM	F1	pass@1	% Win	
CC	57.5	52.0	55.4	56.6	34.5	75.4	55.2
BC	57.2	50.0	49.3	52.0	30.6	64.8	50.7
BC IL	52.6	41.0	46.9	45.4	-	-	-

Table 13: **Base model ablations.** Models are only trained for 100 steps and with other sub-optimal settings, such as the Zephyr format, that were rectified through later ablations.

Task (→)	CLF	Clust.	PairCLF	Rerank	Retrieval	STS	Summ.	Avg.
Metric (→)	Acc.	V-Meas.	AP	MAP	nDCG	Spear.	Spear.	
Dataset # (→)	12	11	3	4	15	10	1	56
Mistral 7B	70.6	43.7	74.0	54.8	35.3	72.9	31.2	54.6
Llama 2 7B	68.1	38.0	64.1	50.2	24.2	67.7	30.5	48.2
GPT-J 6B	70.7	41.4	69.6	53.9	29.7	70.4	29.8	51.9

Dataset (→)	MMLU	GSM8K	BBH	TyDi QA	HumanEval	Avg.
Setup (→)	0 FS	8 FS, CoT	3 FS, CoT	1 FS, GP	0 FS	
Metric (→)	EM	EM	EM	F1	pass@1	
Mistral 7B	35.0	11.0	31.6	20.5	13.8	22.4
Llama 2 7B	35.8	7.0	27.2	21.0	12.9	20.8
GPT-J 6B	27.5	3.5	22.2	8.7	8.0	14.0

Table 14: **Embedding-only models embedding dataset ablations.** NNI = No Natural Instructions, corresponding to not including natural instructions in the data. II = evaluating with the Instructor-XL instructions (Su et al., 2023). Other models use our new structure with domain, intent, and unit depicted in Figure 3. Thus, MEDI2 NNI II and MEDI2 NNI are the same model and only differ in the evaluation instruction set.

Task (→)	CLF	Clust.	PairCLF	Rerank	Retrieval	STS	Summ.	Avg.
Metric (→)	Acc.	V-Meas.	AP	MAP	nDCG	Spear.	Spear.	
Dataset # (→)	12	11	3	4	15	10	1	56
MEDI II	77.1	44.0	83.3	57.0	43.2	79.6	29.4	60.0
MEDI2 NNI II	74.0	43.5	80.5	56.6	46.1	78.4	29.5	59.6
MEDI2 NNI	74.2	44.5	80.7	57.3	49.5	79.6	30.8	61.1
MEDI2	75.1	43.8	80.6	57.5	50.2	81.7	31.9	61.7
MEDI2 + Weights	74.4	42.7	78.4	57.7	50.2	81.4	30.5	61.2

Table 15: **Unified models embedding dataset ablations.** The sequence of Cs and Bs refers to the attention mechanism for (from left to right): Emb instruction, Emb sample, where C=Causal, B=Bidirectional, and Emb=Embedding. WM and M refer to position-weighted mean and mean pooling, respectively. MEDI2BGE corresponds to our MEDI2 dataset with negatives coming from the BGE training dataset MTP (Xiao et al., 2023).

Task (→)	CLF	Clust.	PairCLF	Rerank	Retrieval	STS	Summ.	Avg.
Metric (→)	Acc.	V-Meas.	AP	MAP	nDCG	Spear.	Spear.	
Dataset # (→)	12	11	3	4	15	10	1	56
CCCC WM MEDI	77.9	47.9	81.5	59.0	49.4	80.3	29.4	62.8
CCCC WM MEDI2	76.5	47.0	82.5	59.4	51.4	81.9	30.2	63.2
BBCC M MEDI	79.1	48.8	86.4	59.6	50.3	81.3	31.0	64.0
BBCC M MEDI2	77.0	48.7	86.0	61.0	53.6	83.0	29.1	64.7
BBCC M MEDI2BGE	77.0	48.9	86.9	61.3	53.1	82.8	29.4	64.7
BBCC M E5	79.7	49.5	86.2	59.6	55.3	83.6	29.9	66.0
Dataset (→)	MMLU	GSM8K	BBH	TyDi QA	HumanEval	Alpaca		Avg.
Setup (→)	0 FS	8 FS, CoT	3 FS, CoT	1 FS, GP	0 FS	0 FS, 1.0		
Metric (→)	EM	EM	EM	F1	pass@1	% Win		
CCCC WM MEDI	57.5	45.0	53.1	56.0	32.3	72.9		52.8
CCCC WM MEDI2	57.1	49.0	53.3	55.3	32.3	73.6		53.4
BBCC M MEDI	57.0	46.5	54.5	55.0	30.4	73.8		52.9
BBCC M MEDI2	57.0	50.5	53.8	54.7	32.3	74.7		53.8
BBCC M MEDI2BGE	57.4	48.0	54.7	55.1	32.0	74.7		53.7
BBCC M E5	57.3	47.5	54.2	54.6	33.6	75.4		53.8

Table 16: **Generative dataset ablations.** EP = number of epochs.

Dataset (→)	MMLU	GSM8K	BBH	TyDi QA	HumanEval	Alpaca		Avg.
Setup (→)	0 FS	8 FS, CoT	3 FS, CoT	1 FS, GP	0 FS	0 FS, 1.0		
Metric (→)	EM	EM	EM	F1	pass@1	% Win		
Tülu 2 1 EP	57.5	52.0	55.4	56.6	34.5	75.4		55.2
Tülu 2 2 EP	58.2	53.0	51.9	54.1	37.4	80.5		55.9
OASST 1 EP	53.8	24.0	41.1	28.2	27.4	51.7		37.7
OASST 2 EP	52.4	17.5	45.7	29.2	19.8	61.3		37.7
UltraChat	56.1	43.0	53.8	35.0	25.9	70.3		47.4

Table 17: **Embedding Head.** “→ 1024” refers to down-projecting the final hidden state with a linear layer from 4096 to 1024 dimensions only for embedding tasks.

Task (→)	CLF	Clust.	PairCLF	Rerank	Retrieval	STS	Summ.	Avg.
Metric (→)	Acc.	V-Meas.	AP	MAP	nDCG	Spear.	Spear.	
Dataset # (→)	12	11	3	4	15	10	1	56
No head	77.7	47.9	81.3	58.6	49.2	80.4	29.5	62.7
→ 1024	76.9	47.6	82.1	58.6	48.0	80.1	29.8	62.1
Dataset (→)	MMLU	GSM8K	BBH	TyDi QA	HumanEval	Alpaca		Avg.
Setup (→)	0 FS	8 FS, CoT	3 FS, CoT	1 FS, GP	0 FS	0 FS, 1.0		
Metric (→)	EM	EM	EM	F1	pass@1	% Win		
No head	54.2	42.5	50.6	53.9	28.4	65.5		49.2
→ 1024	53.6	37.0	48.8	54.4	26.6	67.3		48.0

Table 18: **Embedding batch size ablations.** 256 and 4096 indicate the respective embedding batch size. The generative batch size is always 256.

Task ( $\rightarrow$ )	CLF	Clust.	PairCLF	Rerank	Retrieval	STS	Summ.	Avg.
Metric ( $\rightarrow$ )	Acc.	V-Meas.	AP	MAP	nDCG	Spear.	Spear.	
Dataset # ( $\rightarrow$ )	12	11	3	4	15	10	1	56
MEDI2 256	76.5	47.0	82.5	59.4	51.4	81.9	30.2	63.2
MEDI2 4096	77.1	48.0	84.1	60.2	52.8	82.8	30.5	64.2
Dataset ( $\rightarrow$ )	MMLU	GSM8K	BBH	TyDi QA	HumanEval	Alpaca	Avg.	
Setup ( $\rightarrow$ )	0 FS	8 FS, CoT	3 FS, CoT	1 FS, GP	0 FS	0 FS, 1.0		
Metric ( $\rightarrow$ )	EM	EM	EM	F1	pass@1	% Win		
MEDI2 256	57.1	49.0	53.3	55.3	32.3	73.6	53.4	
MEDI2 4096	57.7	48.0	53.2	54.5	32.0	74.3	53.3	

Table 19: **Precision ablations.** BF16 refers to bfloat16 mixed precision and FP32 to float32 precision.

Task ( $\rightarrow$ )	CLF	Clust.	PairCLF	Rerank	Retrieval	STS	Summ.	Avg.
Metric ( $\rightarrow$ )	Acc.	V-Meas.	AP	MAP	nDCG	Spear.	Spear.	
Dataset # ( $\rightarrow$ )	12	11	3	4	15	10	1	56
BF16	79.7	50.2	87.6	60.2	56.5	83.4	30.8	66.5
FP32	79.6	50.3	87.2	59.9	56.1	83.3	30.9	66.3
Dataset ( $\rightarrow$ )	MMLU	GSM8K	BBH	TyDi QA	HumanEval	Alpaca	Avg.	
Setup ( $\rightarrow$ )	0 FS	8 FS, CoT	3 FS, CoT	1 FS, GP	0 FS	0 FS, 1.0		
Metric ( $\rightarrow$ )	EM	EM	EM	F1	pass@1	% Win		
BF16	58.2	51.5	52.8	55.9	37.3	74.4	55.0	
FP32	55.9	52.0	49.9	53.9	31.2	71.3	52.4	

Table 20: **In-batch negatives ablations.**

Task ( $\rightarrow$ )	CLF	Clust.	PairCLF	Rerank	Retrieval	STS	Summ.	Avg.
Metric ( $\rightarrow$ )	Acc.	V-Meas.	AP	MAP	nDCG	Spear.	Spear.	
Dataset # ( $\rightarrow$ )	12	11	3	4	15	10	1	56
Any dataset	79.7	49.8	85.5	59.8	54.9	83.9	30.5	66.0
Same dataset	79.5	48.9	87.4	59.0	56.2	83.0	30.5	66.0
Dataset ( $\rightarrow$ )	MMLU	GSM8K	BBH	TyDi QA	HumanEval	Alpaca	Avg.	
Setup ( $\rightarrow$ )	0 FS	8 FS, CoT	3 FS, CoT	1 FS, GP	0 FS	0 FS, 1.0		
Metric ( $\rightarrow$ )	EM	EM	EM	F1	pass@1	% Win		
Any dataset	56.1	43.5	53.1	46.6	33.5	72.3	50.9	
Same dataset	55.0	45.0	54.4	49.3	29.6	73.4	51.1	

Table 21: **Generative format ablations.**

Dataset ( $\rightarrow$ )	MMLU	GSM8K	BBH	TyDi QA	HumanEval	Alpaca	Avg.	
Setup ( $\rightarrow$ )	0 FS	8 FS, CoT	3 FS, CoT	1 FS, GP	0 FS	0 FS, 1.0		
Metric ( $\rightarrow$ )	EM	EM	EM	F1	pass@1	% Win		
Tulu 2 format	57.5	52.0	55.4	56.6	34.5	75.4	55.2	
Zephyr $\beta$ format	57.3	53.5	52.7	59.1	0.0	71.2	49.0	

Table 22: **Unified models max tokens ablations.** X:Y refers to “maximum tokens allowed for embedding documents during training”: “maximum tokens allowed for queries and documents during embedding evaluation”. The sequence of Cs and Bs refers to the attention mechanism for (from left to right): Emb instruction, Emb sample, where C=Causal, B=Bidirectional, and Emb=Embedding.

Task (→)	CLF	Clust.	PairCLF	Rerank	Retrieval	STS	Summ.	Avg.
Metric (→)	Acc.	V-Meas.	AP	MAP	nDCG	Spear.	Spear.	
Dataset # (→)	12	11	3	4	15	10	1	56
MEDI 2048:512	77.9	47.9	81.5	59.0	49.4	80.3	29.4	62.8
MEDI 2048:4096	77.9	47.9	81.5	59.0	49.4	80.2	31.3	62.8
MEDI 4096:512	76.7	47.3	79.8	58.8	47.0	78.5	30.0	61.3
MEDI 4096:4096	76.8	47.2	79.8	58.8	46.9	78.2	29.9	61.3
MEDI2 BBCC 2048:512	77.0	48.7	86.0	61.0	53.6	83.0	29.1	64.7
MEDI2 BBCC 512:512	76.9	47.6	85.5	61.0	52.8	82.3	28.8	64.1
Dataset (→)	MMLU	GSM8K	BBH	TyDi QA	HumanEval	Alpaca		Avg.
Setup (→)	0 FS	8 FS, CoT	3 FS, CoT	1 FS, GP	0 FS	0 FS, 1.0		
Metric (→)	EM	EM	EM	F1	pass@1	% Win		
MEDI 2048:512/4096	57.4	45.0	53.1	56.0	32.3	72.9		52.8
MEDI 4096:512/4096	53.8	43.0	52.7	54.8	30.1	-		-
MEDI2 BBCC 2048:512	57.0	50.5	53.8	54.7	32.3	74.7		53.8
MEDI2 BBCC 512:512	56.9	46.5	53.1	52.6	31.2	72.8		52.2

Table 23: **Loss ablations.** E.g. Mix (32 → 8) corresponds to token level loss across 32 samples and then sample level loss across 8 sub-batches for a total batch size of 256. E.g. 2.4 refers to the loss ratio of the 1st step:  $\mathcal{L}_{Emb}/\mathcal{L}_{Gen}$ .

Task (→)	CLF	Clust.	PairCLF	Rerank	Retrieval	STS	Summ.	Avg.
Metric (→)	Acc.	V-Meas.	AP	MAP	nDCG	Spear.	Spear.	
Dataset # (→)	12	11	3	4	15	10	1	56
E5S Token 2.4	79.5	50.1	86.5	60.0	55.6	83.2	30.3	66.1
E5S Token 6.0	79.7	50.2	87.6	60.2	56.5	83.4	30.8	66.5
E5S Mix (32 → 8) 4.1	79.4	50.5	87.2	60.5	57.4	83.4	30.4	66.7
Dataset (→)	MMLU	GSM8K	BBH	TyDi QA	HumanEval	Alpaca		Avg.
Setup (→)	0 FS	8 FS, CoT	3 FS, CoT	1 FS, GP	0 FS	0 FS, 1.0		
Metric (→)	EM	EM	EM	F1	pass@1	% Win		
E5S Token 2.4	57.9	48.5	53.5	56.5	35.2	75.0		54.4
E5S Token 6.0	58.2	51.5	52.8	55.9	37.3	74.4		55.0
E5S Mix (32 → 8) 4.1	57.6	57.0	54.8	55.4	32.8	74.8		55.4
MEDI2 Mix (4 → 64) 11.7	57.0	48.0	53.7	55.0	35.8	67.6		52.9
MEDI2 Mix (32 → 8) 10.2	57.0	50.5	53.8	54.7	32.3	74.7		53.8

## J GRITLM MTEB FULL RESULTS

Table 24: MTEB full results from Table 1.

Dataset	Gen-only	Emb-only	GRITLM	
			7B	8x7B
AmazonCounterfactualClassification	70.06	82.55	81.18	80.48
AmazonPolarityClassification	74.74	96.19	96.52	96.32
AmazonReviewsClassification	38.63	57.28	57.81	57.18
Banking77Classification	71.25	88.73	88.47	87.46
EmotionClassification	36.61	51.83	52.81	50.06
ImdbClassification	73.94	94.58	95.00	94.32
MassiveIntentClassification	66.82	79.37	80.78	79.72
MassiveScenarioClassification	71.27	81.20	82.09	81.09
MTOPDomainClassification	85.40	96.72	96.16	95.29
MTOPIntentClassification	75.60	87.19	87.13	87.08
ToxicConversationsClassification	66.36	68.37	70.80	70.89
TweetSentimentExtractionClassification	54.61	61.91	64.78	62.48
ArxivClusteringP2P	45.40	50.87	51.67	50.72
ArxivClusteringS2S	29.86	47.35	48.11	48.01
BiorxivClusteringP2P	33.45	40.18	40.87	41.41
BiorxivClusteringS2S	23.02	39.60	39.80	38.67
MedrxivClusteringP2P	27.49	36.61	36.52	36.54
MedrxivClusteringS2S	23.17	37.28	36.80	37.24
RedditClustering	23.28	63.52	61.30	63.01
RedditClusteringP2P	55.00	67.81	67.26	65.86
StackExchangeClustering	47.14	75.53	77.33	74.41
StackExchangeClusteringP2P	33.95	46.22	41.33	38.52
TwentyNewsgroupsClustering	18.15	56.8	55.70	57.16
SprintDuplicateQuestions	51.57	93.37	93.00	91.24
TwitterSemEval2015	50.60	80.61	81.08	77.21
TwitterURLCorpus	60.36	87.20	87.40	86.45
AskUbuntuDupQuestions	49.02	68.13	67.34	65.60
MindSmallReranking	27.83	32.19	31.81	32.84
SciDocsRR	56.65	87.00	86.84	86.43
StackOverflowDupQuestions	38.42	55.48	55.96	54.33
ArguAna	35.96	62.95	63.24	59.49
ClimateFEVER	8.96	31.09	30.91	28.69
CQADupstackRetrieval	7.20	50.83	49.42	47.63
DBPedia	2.15	47.06	46.60	46.54
FEVER	5.02	85.41	82.74	85.02
FiQA2018	6.27	60.22	59.95	49.89
HotpotQA	6.67	79.15	79.40	73.83
MSMARCO	0.66	41.55	41.96	35.55
NFCorpus	3.74	41.69	40.89	39.05
NQ	2.14	69.46	70.30	63.87
QuoraRetrieval	64.42	89.08	89.47	87.70
SCIDOCS	2.32	24.86	24.41	23.06
SciFact	35.58	78.92	79.17	77.02
Touche2020	3.06	24.30	27.93	27.97
TRECCOVID	20.92	75.29	74.8	81.07
BIOSSES	70.87	86.20	86.35	87.34
SICK-R	58.95	83.03	83.13	80.56
STS12	44.25	78.07	77.34	73.69
STS13	64.22	85.98	85.04	85.82

2052	STS14	52.24	83.92	82.91	82.05
2053	STS15	64.53	89.18	88.13	88.8
2054	STS16	65.89	86.83	86.24	86.2
2055	STS17	69.64	89.7	90.13	91.46
2056	STS22	57.29	68.41	68.63	69.21
2057	STSBenchmark	53.89	86.74	85.64	87.43
2058	SummEval	21.14	30.18	30.37	29.82
2059	<b>Average</b>	<b>41.21</b>	<b>66.82</b>	<b>66.76</b>	<b>65.66</b>

### K REDUCING EMBEDDING TRAINING MEMORY

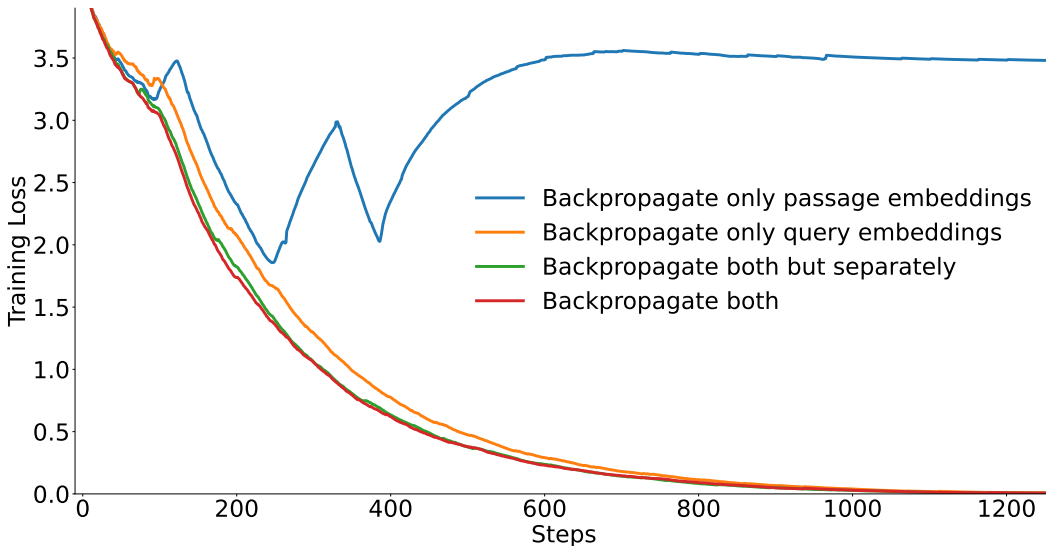


Figure 9: **Embedding memory ablations.** Passage corresponds to both positive and document embeddings. Loss is smoothed with exponential moving average smoothing and a weight of 0.99.

Generative training only requires sufficient memory to perform a forward and backward pass on a single training sample of a given sequence length. Meanwhile, naive embedding training with in-batch negatives requires sufficient memory to accommodate a forward and a backward pass on  $3 * bs$  samples. The 3 corresponds to the need for passing a triplet of a query, a positive, and a negative document (Equation 1). The batch size ( $bs$ ) factor corresponds to the need for forwarding all samples together as regular gradient accumulation does not work with in-batch negatives. Below we outline the strategies we employ to reduce these memory needs.

**Triplet** As the full triplet is only required for loss calculation (Equation 1), it can be split across separate forward and backward passes. To avoid the memory requirements of gradients in PyTorch Autograd (Paszke et al., 2019), this requires two additional forward passes without gradients. Simplified code representing this procedure is depicted in Listing 1. In our training, it was sufficient to only split the triplet into two parts: query and passages, where passages consist of both a positive and a negative document. Thus, we only incur the cost of one additional forward pass without gradients on the query. Alternatively, one could only backpropagate on a subset of the embeddings, however, we show in Figure 9 that this leads to worse performance.

**In-batch negatives** There are two strategies to reduce the batch size memory requirement to that of a single batch while using nearly unlimited in-batch negatives. **(1) Distributed Training:** The best strategy is to distribute the training across up to  $bs$  GPUs. The representations can then be gathered across GPUs to compute the contrastive loss with in-batch negatives. **(2) GradCache:** If enough GPUs are not available, GradCache (Gao et al., 2021b) can be used. GradCache maintains

Listing 1: **Splitting of the embedding pass to save memory, simplified.**

```

2106
2107
2108
2109 def distributed_contrastive_loss(q, p, n):
2110     # Gather in-batch negatives across devices...
2111     # Compute contrastive loss...
2112
2113     # Split triplet into three forward passes
2114     pos_rep = model(pos)
2115     with torch.no_grad():
2116         q_rep = model(query)
2117         neg_rep = model(neg)
2118
2119     # Only perform backward pass on positive documents
2120     loss = distributed_contrastive_loss(q_rep, pos_rep, neg_rep)
2121     loss.backward()
2122
2123     pos_rep = pos_rep.detach()
2124     # Perform forward + backward on negatives & reuse rest
2125     neg_rep = model(neg)
2126     loss = distributed_contrastive_loss(q_rep, pos_rep, neg_rep)
2127     loss.backward()
2128
2129     # Perform forward + backward on queries & reuse rest
2130     neg_rep = neg_rep.detach()
2131     q_rep = model(query)
2132     loss = distributed_contrastive_loss(q_rep, pos_rep, neg_rep)
2133     loss.backward()

```

in-batch negatives while allowing computation of gradients for each triplet at a time, thus effectively corresponding to gradient accumulation for contrastive loss. However, it comes at the cost of additional forward passes.

Across training runs, we make use of all three strategies (splitting, distributed training, GradCache).

## L HYPERPARAMETERS

We finetune all parameters of our models for up to 1253 steps. Our learning rate is  $2e-5$ , we use 3% of steps for linear warm-up of the learning rate and decay it linearly to 0 over training. To save memory, we use PyTorch FSDP (Zhao et al., 2023), gradient checkpointing, BF16 mixed precision training, and strategies outlined in Appendix K. During training, we use a sequence length of 2048 for generative samples, 256 for embedding queries, and 2048 for embedding documents unless otherwise specified. We finetune using the Adam optimizer (Kingma & Ba, 2017) with  $\beta_1=0.9$  and  $\beta_2=0.999$  and no weight decay. We also use Flash-Attention 2 (Dao et al., 2022; Dao, 2023) via PyTorch SDPA.

We evaluate models using the settings put forth by the creators of MTEB (Muennighoff et al., 2023c), Tülu (Iverson et al., 2023; Wang et al., 2024) and HumanEvalSynthesize (Muennighoff et al., 2023a; Zhuo et al., 2024). For MTEB, we evaluate using a maximum sequence length of 512 unless otherwise specified.

## M EMBEDDING INSTRUCTION FOR GENERATIVE MODELS

As prior instruction-tuned models have been trained without an embedding objective, it is unclear whether one should add an instruction when evaluating them on embedding tasks. We benchmark the Mistral 7B instruct model on MTEB with and without instruction in Table 25. We find that performance is around the same, however, adding instructions performs slightly better. Thus, we add an instruction for all instruction-tuned models when benchmarking their embedding performance.



Table 25: **Benchmarking the benefit of an embedding instruction for generative instruction-tuned models.** When an instruction is used (“Mistral Instruct w/”), we use the default instructions from Instructor XL with the prompt template of the Mistral Instruct model. For no instruction (“Mistral Instruct w/o”), the procedure is the same as for the base model (“Mistral”)

Task (→)	CLF	Clust.	PairCLF	Rerank	Retrieval	STS	Summ.	Avg.
Metric (→)	Acc.	V-Meas.	AP	MAP	nDCG	Spear.	Spear.	
Dataset # (→)	12	11	3	4	15	10	1	56
Mistral	63.5	34.6	53.5	43.2	13.2	57.4	19.7	40.5
Mistral Instruct w/o	65.4	35.6	60.2	44.6	16.8	61.1	25.9	43.3
Mistral Instruct w/	67.1	34.6	59.6	44.8	16.3	63.4	25.9	43.7

## N HUMAN EVAL FORMAT

Table 26: **HumanEvalSynthesize with different formats using Tülu 2 7B.**

In Tülu 2 (Iverson et al., 2023), models are evaluated on HumanEval (Chen et al., 2021) without the model’s chat format. As this does not reflect the intended usage of the models, we instead use the appropriate chat format for evaluating HumanEval. To do so, we use the instructions and evaluation procedure from HumanEval-Synthesize (Muennighoff et al., 2023a). In Table 26 we benchmark the impact this has on performance for the Tülu 2 7B model (Iverson et al., 2023). We find that the performance is around equivalent and thus use the chat format for all evaluations of chat models. For non-chat models, we use the original HumanEval continuation format as proposed by Chen et al. (2021)

Format	Tülu 2 7B	
	No Chat	Chat
Pass@1	23.4	24.5
Pass@10	32.4	31.3

## O EMBEDDING IN FP32 VS BF16

We perform all training and evaluations in BF16 (bfloat16) mixed precision to speed up computations. We verified that it performs comparably to FP32 (float32) on MTEB in Table 27. Note that pooling and subsequent similarity computations are still in FP32.

Table 27: **Embeddings in FP32 vs BF16.** Benchmarking of the raw Mistral 7B model. “FP32” corresponds to doing all computations in float32 precision. “BF16” and “BF16 Cache” corresponds to doing most operations in bfloat16 except for operations that PyTorch auto casts to float32 (e.g. normalization), pooling and similarity computations. For “BF16 Cache”, we cast the embeddings after pooling to BF16 and then back to FP32 before similarity computations. This corresponds to locally caching the embeddings in BF16 to save storage and then casting them to FP32 at inference.

Task (→)	CLF	Clust.	PairCLF	Rerank	Retrieval	STS	Summ.	Avg.
Metric (→)	Acc.	V-Meas.	AP	MAP	nDCG	Spear.	Spear.	
Dataset # (→)	12	11	3	4	15	10	1	56
FP32	63.46	34.62	53.56	43.24	13.26	57.38	19.87	40.51
BF16	63.47	34.60	53.52	43.24	13.24	57.38	19.68	40.50
BF16 Cache	63.47	34.56	53.52	43.25	13.11	57.38	19.71	40.46

## P UNRELIABILITY OF MT-BENCH

We experiment with using MT-Bench with its recommended absolute scores for our generative evaluation (Zheng et al., 2023). However, we find that as soon as we switch the LLM Evaluator from GPT-4 to GPT-4 Turbo, the scores change significantly (Table 28). GPT-4 is a closed-source model with changes happening behind the scenes that users may not know about (Chen et al., 2023). Thus, if OpenAI decides to change GPT-4, all existing MT-Bench absolute scores would essentially become obsolete. The same applies if the API is retired. To alleviate this, we also experiment with using Zephyr 7B  $\beta$  (Tunstall et al., 2023) and Llama 2 70B Chat (Touvron et al., 2023) as evaluators, however, we find them to often not provide any rating as they struggle to understand the prompt. While AlpacaEval (Dubois et al., 2023; Li et al., 2023b), which we use, shares some of these problems, its comparison-based evaluation is more stable. This is because comparing if one generation is better than another generally has an objective ground truth solution. Meanwhile, there is no objective solution as to whether an absolute score of a given generation should be 3 or 4 (MT-Bench has eleven levels from 0-10). This is up to the subjective value system of the evaluator.

Table 28: **Using GPT-4 vs GPT-4 Turbo as a judge for MT-Bench.** Each evaluator is provided with the same generations of the same instruction-tuned model.

	GPT-4	GPT-4 Turbo	Drop
Turn 1	4.08	3.05	25%
Turn 2	2.64	1.88	29%
Avg.	3.36	2.48	26%

## Q LIMITATIONS AND FUTURE WORK

**Efficiency** As mentioned in §1, training using GRIT requires more compute than only embedding or only generative training as two forward and backward passes are required. As finetuning is generally cheaper than pretraining, this is not a major problem, but efficiency improvements would nonetheless be worthwhile. One potential way to improve efficiency would be to extract the embedding and generative signal from the same samples, rather than separate samples. This could halve the number of forward passes required, yet due to the different loss functions, it may not make the backward passes significantly faster.

**Performance improvements** While we find that GRITLM performs strongly on embedding and generative tasks (§3.2), there have been many recent models with even stronger performance in either embedding or generative tasks; yet not the combination of both. A natural future work would therefore be extending the GRIT approach to more recent models, such as the Llama-3 series of models (Dubey et al., 2024) to build stronger models that can handle both embedding and generation.

**Caching improvements** As we outline in §5, the caching variants with GRITLM suffer from attention mismatch problems. Further, doc caching requires a significant amount of extra storage. While storage is usually cheap, it may nonetheless be prohibitively expensive for very large indices. One promising avenue for future work is improving caching with GRITLM, such as via finetuning with caching, such that it learns to deal with the mismatch problem.

**GRITLM Agents** Future work may consider using the embedding capability to let the generative model initiate a search over an index when it deems necessary. Currently, this is often accomplished via external retrieval plugins. Such plugins are no longer necessary if the model can retrieve on its own. Teaching the model to invoke its own embedding capability likely requires additional finetuning (just like teaching it to invoke an external plugin (Schick et al., 2023)). A sample could look something like:

```
“<|user|>\nWhat is the capital of Japan?\n<|internal|>\nI am not sure I know this. Let me produce an embedding for it and search for the answer. Retrieve answers for this query.\n<|embed|>\nWhat is the capital of Japan?\n<|output|>\nTokyo, Japan’s busy capital, mixes the ultramodern and the traditional..\n<|assistant|>\n\nThe capital of Japan is Tokyo.\n</s>”
```

**Pretraining** For our experiments we take an off-the-shelf pretrained language model. However, it should also be possible to use the GRIT approach to pretrain from scratch. As labeled embedding data is likely too scarce for pretraining, one could either rely on unsupervised approaches for the embedding objective, such as RetroMAE (Xiao et al., 2022; Xiao & Liu, 2022), or use methods like data augmentation (Dhole et al., 2022), pruning (Xia et al., 2023) or multi-epoch training to deal with the data constraint (Muennighoff et al., 2023b; Luukkonen et al., 2023).

**Format Efficiency** Our format in Figure 3 is inefficient, as encoding the embedding format, `<s><|user|>\n<|embed|>\n`, requires 13 tokens and encoding the generative format, `<s><|user|>\n<|assistant|>\n</s>`, requires 15 tokens. Using special tokens could simplify this and thus make training and inference slightly cheaper.

**Training efficiency: Packing and Reusing** It is common to pack samples during generative instruction tuning to maximize efficiency (Chung et al., 2022; Muennighoff et al., 2023d). Packing embedding samples during training should also be possible by ensuring attention is only paid to each respective sample. Going even further is it possible to pack generative and embedding training data into the same sample and reuse the same sample for both tasks? This could look similar to the example provided in “GRITLM Agents” with the generative loss applied over the assistant response and the contrastive loss applied to the representation of the text following “`<|embed|>`”. By reusing samples it may be possible to significantly decrease the resources needed for GRIT.

## R DATASET COMPOSITION

Table 29: **E5S dataset composition.**

Dataset (↓)	Num samples
DuReader (Qiu et al., 2022)	86,395
ELI5 (Fan et al., 2019)	50293
FEVER (Thorne et al., 2018)	71,257
GPT4 Bitext (Wang et al., 2024)	89,324
GPT4 P2P (Wang et al., 2024)	16,842
GPT4 P2S (Wang et al., 2024)	121,878
GPT4 Retrieval (Wang et al., 2024)	166,602
GPT4 S2S (Wang et al., 2024)	13,481
GPT4 STS (Wang et al., 2024)	98,626
HotpotQA (Yang et al., 2018)	68,659
NLI (Gao et al., 2022)	275,601
MIRACL (Zhang et al., 2022)	40,203
MSMARCO (Bajaj et al., 2018)	244,582
MSMARCO Doc (Bajaj et al., 2018)	71,594
Mr. TyDi (Zhang et al., 2021)	48,729
NQ (Kwiatkowski et al., 2019)	71,408
S2ORC (Lo et al., 2020)	80,000
SQuAD (Rajpurkar et al., 2016)	87,599
T2Ranking (Xie et al., 2023)	112,335
TriviaQA (Karpukhin et al., 2020)	60,296
Quora (DataCanary et al., 2017)	14,926
<b>Total</b>	<b>1,890,630</b>

2322  
2323  
2324  
2325  
2326  
2327  
2328  
2329  
2330  
2331  
2332  
2333  
2334  
2335  
2336  
2337  
2338  
2339  
2340  
2341  
2342  
2343  
2344  
2345  
2346  
2347  
2348  
2349  
2350  
2351  
2352  
2353  
2354  
2355  
2356  
2357  
2358  
2359  
2360  
2361  
2362  
2363  
2364  
2365  
2366  
2367  
2368  
2369  
2370  
2371  
2372  
2373  
2374  
2375

Table 30: **MEDI2 dataset composition.**

MEDI Dataset (↓)	Num samples
AGNews (Zhang et al., 2016)	199,792
Altlex (Hidey & McKeown, 2016)	112,602
Amazon QA (Gupta et al., 2019)	199,180
Amazon Review (Keung et al., 2020)	198,298
CC News (Hamborg et al., 2017)	190,503
CNN/Dailymail (Fabbri et al., 2021)	189,407
COCO Captions (Chen et al., 2015)	82,783
ELI5 (Fan et al., 2019)	196,572
FEVER KILT (Thorne et al., 2018; Petroni et al., 2021)	71,257
Flickr 30k (Young et al., 2014)	31,783
Gigaword (Rush et al., 2015; Graff et al., 2003)	200,000
GooAQ (Khashabi et al., 2021)	199,981
HotpotQA KILT (Yang et al., 2018; Petroni et al., 2021)	65,351
NLI (Gao et al., 2022)	277,195
MSMARCO (Bajaj et al., 2018)	491,980
MedMCQA (Pal et al., 2022)	156,905
Multi-LexSum (Shen et al., 2022)	2,771
NPR (Team, 2021b)	193,399
NQ (Kwiatkowski et al., 2019)	73,226
PAQ (Lewis et al., 2021b)	190,162
PubMedQA (Jin et al., 2019)	190,481
Reddit (Team, 2021c)	196,247
S2ORC (Lo et al., 2020)	193,458
SQuAD (Rajpurkar et al., 2016)	84,105
SciTLDR (Cachola et al., 2020)	1,742
SearchQA (Dunn et al., 2017)	114,520
Sentence Compression (Filippova & Altun, 2013)	179,996
SimpleWiki (Coster & Kauchak, 2011)	102,035
StackExchange (Team, 2021a)	201,050
SuperNI (300 datasets) (Wang et al., 2022c)	2,682,465
SPECTER (Cohan et al., 2020)	684,000
T-REx KILT (ElSahar et al., 2018; Petroni et al., 2021)	191,383
Quora (DataCanary et al., 2017)	101,762
WikiAnswers (Fader et al., 2014)	200,000
WikiHow (Koupae & Wang, 2018)	128,542
XSum (Narayan et al., 2018)	190,427
Yahoo (Zhang et al., 2016)	198,346
Zeroshot KILT (Levy et al., 2017; Petroni et al., 2021)	124,547
<b>Total</b>	<b>9,084,806</b>

2376 S DATASET SAMPLES  
2377

---

2378 **Query instruction:**

---

2379 Represent the sentence for retrieving supporting documents;  
2380

---

2381 **Query sample:**

---

2382 what two plates form the san andreas fault  
2383

---

2384 **Positive instruction:**

---

2385 Represent the document for retrieval;  
2386

---

2387 **Positive sample:**

---

2388 The San Andreas Fault marks the junction between the North American and Pacific Plates. The fault  
2389 is 1300 km long, extends to at least 25 km in depth, and has a north west-south east trend. It is  
2390 classified as a right lateral (dextral) strike-slip fault. Loading the player ...  
2391

---

2392 **Negative instruction:**

---

2393 Represent the document for retrieval;  
2394

---

2395 **Negative sample:**

---

2396 The San Andreas Fault is the sliding boundary between the Pacific Plate and the North American  
2397 Plate. It slices California in two from Cape Mendocino to the Mexican border. San Diego, Los  
2398 Angeles and Big Sur are on the Pacific Plate.  
2399

---

2400 Figure 10: **MEDI sample.**  
2401  
2402  
2403  
2404  
2405  
2406  
2407  
2408  
2409  
2410  
2411  
2412  
2413  
2414  
2415  
2416  
2417  
2418  
2419  
2420  
2421  
2422  
2423  
2424  
2425  
2426  
2427  
2428  
2429

2430  
2431  
2432  
2433  
2434  
2435  
2436  
2437  
2438  
2439  
2440  
2441  
2442  
2443  
2444  
2445  
2446  
2447  
2448  
2449  
2450  
2451  
2452  
2453  
2454  
2455  
2456  
2457  
2458  
2459  
2460  
2461  
2462  
2463  
2464  
2465  
2466  
2467  
2468  
2469  
2470  
2471  
2472  
2473  
2474  
2475  
2476  
2477  
2478  
2479  
2480  
2481  
2482  
2483

---

**Query instruction:**

Represent this question to retrieve a fitting Wikipedia passage (formal)

---

**Query sample:**

which two plates meet along the west coast of the USA

---

**Positive instruction:**

Represent this Wikipedia text in order to get a user query which it answers!

---

**Positive sample:**

on to a transitional deformation zone in the Chersky Range, then the Ulakhan Fault between it and the Okhotsk Plate, and finally the Aleutian Trench to the end of the Queen Charlotte Fault system. The westerly boundary is the Queen Charlotte Fault running offshore along the coast of Alaska and the Cascadia subduction zone to the north, the San Andreas Fault through California, the East Pacific Rise in the Gulf of California, and the Middle America Trench to the south. On its western edge, the Farallon Plate has been subducting

---

**Negative instruction:**

Represent this passage to easily find a natural-written user question that can be answered by it.

---

**Negative sample:**

the continental margin.  
Types.  
There are two types of continental margins: active and passive margins.  
Active margins are typically associated with lithospheric plate boundaries. These active margins can be convergent or transform margins, and are also places of high tectonic activity, including volcanoes and earthquakes. The West Coast of North America and South America are active margins. Active continental margins are typically narrow from coast to shelf break, with steep descents into trenches. Convergent active margins occur where oceanic plates meet continental

---

Figure 11: **MEDI2 sample.**

2484  
2485  
2486  
2487  
2488  
2489  
2490  
2491  
2492  
2493  
2494  
2495  
2496  
2497  
2498  
2499  
2500  
2501  
2502  
2503  
2504  
2505  
2506  
2507  
2508  
2509  
2510  
2511  
2512  
2513  
2514  
2515  
2516  
2517  
2518  
2519  
2520  
2521  
2522  
2523  
2524  
2525  
2526  
2527  
2528  
2529  
2530  
2531  
2532  
2533  
2534  
2535  
2536  
2537

---

**Query instruction:**

Given a question, retrieve Wikipedia passages that answer the question

---

**Query sample:**

which two plates meet along the west coast of the USA

---

**Positive sample:**

North American Plate boundary is the Queen Charlotte Fault running offshore along the coast of Alaska and the Cascadia subduction zone to the north, the San Andreas Fault through California, the East Pacific Rise in the Gulf of California, and the Middle America Trench to the south. On its western edge, the Farallon Plate has been subducting under the North American Plate since the Jurassic Period. The Farallon Plate has almost completely subducted beneath the western portion of the North American Plate leaving that part of the North American Plate in contact with the Pacific Plate as the San Andreas Fault. The Juan

---

**Negative sample:**

Caribbean Plate Caribbean Plate The Caribbean Plate is a mostly oceanic tectonic plate underlying Central America and the Caribbean Sea off the north coast of South America. Roughly 3.2 million square kilometers (1.2 million square miles) in area, the Caribbean Plate borders the North American Plate, the South American Plate, the Nazca Plate and the Cocos Plate. These borders are regions of intense seismic activity, including frequent earthquakes, occasional tsunamis, and volcanic eruptions. The northern boundary with the North American plate is a transform or strike-slip boundary which runs from the border area of Belize, Guatemala (Motagua Fault), and Honduras in Central

---

Figure 12: **E5 sample.** The E5 dataset does not use instructions for documents, thus the positive and negative samples do not have instructions.

---

**Instruction:**

Q: Lloyd, Mark, and Michael have their Pokemon cards collection. Currently, Mark has thrice as many cards as Lloyd but has 10 fewer cards than Michael. If Michael has 100 cards now, how many more cards should they collect so that all three of them will have a total of 300 cards?

A: 80

Explain how we arrive at this answer:

---

**Sample:**

Explanation: Mark has 10 fewer cards than Michael so Mark has 100 cards - 10 cards = 90 cards. So, Lloyd has 90 cards / 3 = 30 cards. All three of them have 90 cards + 30 cards + 100 cards = 220 cards. Thus, they need to collect 300 cards - 220 cards = 80 more cards.

---

Figure 13: **Tülu 2 sample.**

## 2538 T EVALUATION PROMPTS

2539

## 2540 T.1 EMBEDDING PROMPTS

2541

2542 Table 31 contains the prompt for each MTEB dataset when training on the E5 dataset, which are the  
 2543 same instructions as used in Wang et al. (2024). Table 32 contains the MTEB prompts we use when  
 2544 training on MEDI2, which we wrote ourselves. For models trained on MEDI, we use the instructions  
 2545 for Instructor-XL from Su et al. (2023).

2546

2547 Table 31: **Instructions used for evaluation on the MTEB benchmark when training with the E5**  
 2548 **dataset.** “STS\*” indicates we use the same instructions for all the STS tasks. For retrieval datasets,  
 2549 we do not use an instruction for the document and only display the query instruction.

2550

2551 Task Name	2552 Instruction
2553 AmazonCounterfactualClassif.	2554 Classify a given Amazon customer review text as either counter- 2555 factual or not-counterfactual
2556 AmazonPolarityClassification	2557 Classify Amazon reviews into positive or negative sentiment
2558 AmazonReviewsClassification	2559 Classify the given Amazon review into its appropriate rating 2560 category
2561 Banking77Classification	2562 Given a online banking query, find the corresponding intents
2563 EmotionClassification	2564 Classify the emotion expressed in the given Twitter message 2565 into one of the six emotions: anger, fear, joy, love, sadness, and 2566 surprise
2567 ImdbClassification	2568 Classify the sentiment expressed in the given movie review text 2569 from the IMDB dataset
2570 MassiveIntentClassification	2571 Given a user utterance as query, find the user intents
2572 MassiveScenarioClassification	2573 Given a user utterance as query, find the user scenarios
2574 MTOPDomainClassification	2575 Classify the intent domain of the given utterance in task-oriented 2576 conversation
2577 MTOPIntentClassification	2578 Classify the intent of the given utterance in task-oriented conver- 2579 sation
2580 ToxicConversationsClassif.	2581 Classify the given comments as either toxic or not toxic
2582 TweetSentimentClassification	2583 Classify the sentiment of a given tweet as either positive, negative, 2584 or neutral
2585 ArxivClusteringP2P	2586 Identify the main and secondary category of Arxiv papers based 2587 on the titles and abstracts
2588 ArxivClusteringS2S	2589 Identify the main and secondary category of Arxiv papers based 2590 on the titles
2591 BiorxivClusteringP2P	Identify the main category of Biorxiv papers based on the titles and abstracts
BiorxivClusteringS2S	Identify the main category of Biorxiv papers based on the titles
MedrxivClusteringP2P	Identify the main category of Medrxiv papers based on the titles and abstracts
MedrxivClusteringS2S	Identify the main category of Medrxiv papers based on the titles
RedditClustering	Identify the topic or theme of Reddit posts based on the titles
RedditClusteringP2P	Identify the topic or theme of Reddit posts based on the titles and posts



2592	StackExchangeClustering	Identify the topic or theme of StackExchange posts based on the titles
2593		
2594	StackExchangeClusteringP2P	Identify the topic or theme of StackExchange posts based on the given paragraphs
2595		
2596		
2597	TwentyNewsgroupsClustering	Identify the topic or theme of the given news articles
2598		
2599	SprintDuplicateQuestions	Retrieve duplicate questions from Sprint forum
2600	TwitterSemEval2015	Retrieve tweets that are semantically similar to the given tweet
2601	TwitterURLCorpus	Retrieve tweets that are semantically similar to the given tweet
2602		
2603	AskUbuntuDupQuestions	Retrieve duplicate questions from AskUbuntu forum
2604	MindSmallReranking	Retrieve relevant news articles based on user browsing history
2605	SciDocsRR	Given a title of a scientific paper, retrieve the titles of other relevant papers
2606		
2607		
2608	StackOverflowDupQuestions	Retrieve duplicate questions from StackOverflow forum
2609		
2610	ArguAna	Given a claim, find documents that refute the claim
2611	ClimateFEVER	Given a claim about climate change, retrieve documents that support or refute the claim
2612		
2613	CQADupstackRetrieval	Given a question, retrieve detailed question descriptions from Stackexchange that are duplicates to the given question
2614		
2615	DBPedia	Given a query, retrieve relevant entity descriptions from DBPedia
2616	FEVER	Given a claim, retrieve documents that support or refute the claim
2617	FiQA2018	Given a financial question, retrieve user replies that best answer the question
2618		
2619		
2620	HotpotQA	Given a multi-hop question, retrieve documents that can help answer the question
2621		
2622	MSMARCO	Given a web search query, retrieve relevant passages that answer the query
2623		
2624	NFCorpus	Given a question, retrieve relevant documents that best answer the question
2625		
2626		
2627	NQ	Given a question, retrieve Wikipedia passages that answer the question
2628		
2629	QuoraRetrieval	Given a question, retrieve questions that are semantically equivalent to the given question
2630		
2631	SCIDOCS	Given a scientific paper title, retrieve paper abstracts that are cited by the given paper
2632		
2633	SciFact	Given a scientific claim, retrieve documents that support or refute the claim
2634		
2635	Touche2020	Given a question, retrieve detailed and persuasive arguments that answer the question
2636		
2637	TRECCOVID	Given a query on COVID-19, retrieve documents that answer the query
2638		
2639		
2640		
2641	STS*	Retrieve semantically similar text.
2642		
2643	SummEval	Given a news summary, retrieve other semantically similar summaries
2644		
2645		

2646  
2647  
2648  
2649  
2650  
2651  
2652  
2653  
2654  
2655  
2656  
2657  
2658  
2659  
2660  
2661  
2662  
2663  
2664  
2665  
2666  
2667  
2668  
2669  
2670  
2671  
2672  
2673  
2674  
2675  
2676  
2677  
2678  
2679  
2680  
2681  
2682  
2683  
2684  
2685  
2686  
2687  
2688  
2689  
2690  
2691  
2692  
2693  
2694  
2695  
2696  
2697  
2698  
2699

**Table 32: Instructions used for evaluation on the MTEB benchmark when training with the MEDI2 dataset.** For asymmetric datasets, Q refers to instructions for queries, while D refers to document instructions.

Task Name	Instruction
AmazonCounterfactualClassification	Represent the text to find another sentence with the same counterfactuality, e.g. sentences with "would", "wish", etc. should match with other sentences of that kind.
AmazonPolarityClassification	Represent the review for finding another Amazon review with the same sentiment (positive / negative)
AmazonReviewsClassification	Represent the review for finding another Amazon review with the same rating
Banking77Classification	Represent the text for finding another one-sentence banking query with the same intent
EmotionClassification	Represent the text for finding another one-sentence text with the same emotion
ImdbClassification	Represent the text for finding another one-sentence movie review with the same sentiment
MassiveIntentClassification	Represent the text for finding another text of a few words with the same intent
MassiveScenarioClassification	Represent the text for finding another text of a few words about the same scenario
MTOPDomainClassification	Represent the text for finding another text of a few words about the same domain
MTOPIntentClassification	Represent the text for finding another text of a few words with the same intent
ToxicConversationsClassification	Represent the text for finding another comment of up to a passage in length with the same level of toxicity (either toxic or not toxic)
TweetSentimentExtractionClassification	Represent the tweet for finding another tweet with the same sentiment (positive / neutral / negative)
ArxivClusteringP2P	Represent the text to find another arXiv title with abstract (concatenated) about the same topic
ArxivClusteringS2S	Represent the text to find another arXiv title about the same topic
BiorxivClusteringP2P	Represent the text to find another bioRxiv title with abstract (concatenated) about the same topic
BiorxivClusteringS2S	Represent the text to find another bioRxiv title about the same topic
MedrxivClusteringS2S	Represent the text to find another medRxiv title about the same topic
MedrxivClusteringP2P	Represent the text to find another medRxiv title with abstract (concatenated) about the same topic
RedditClustering	Represent the text to find another Reddit community title that stems from the same subreddit
RedditClusteringP2P	Represent the text to find another Reddit community title with post (concatenated) from the same subreddit
StackExchangeClustering	Represent the text to find another StackExchange title that stems from the same StackExchange

2700	StackExchangeClusteringP2P	Represent the text to find another StackExchange title with post (concatenated) that stems from the same StackExchange
2701		
2702	TwentyNewsgroupsClustering	Represent the title to find a similar news title from the same newsgroup
2703		
2704		
2705	SprintDuplicateQuestions	Represent the question to be matched with another duplicate user question from the Sprint community forum
2706		
2707	TwitterSemEval2015	Represent the tweet to find another tweet that is a paraphrase of it
2708		
2709	TwitterURLCorpus	Represent the tweet to find another tweet that is a paraphrase of it
2710		
2711	ArguAna Q	Represent the passage to find a passage with a counter-argument about the same topic to it
2712		
2713	ArguAna D	Represent the passage to find a passage with a counter-argument about the same topic to it
2714		
2715	ClimateFEVER Q	Represent the climate-based claim to find a Wikipedia abstract to support it
2716		
2717	ClimateFEVER D	Represent the Wikipedia abstract to find a climate-related claim that it supports
2718		
2719		
2720	CQADupstackAndroidRetrieval Q	Represent the title of a user question to find a duplicate user question title with body from the Android StackExchange forum
2721		
2722	CQADupstackAndroidRetrieval D	Represent the question title with body posted by a user to find a duplicate user question title from the Android StackExchange forum
2723		
2724		
2725	CQADupstackEnglishRetrieval Q	Represent the title of a user question to find a duplicate user question title with body from the English StackExchange forum
2726		
2727	CQADupstackEnglishRetrieval D	Represent the question title with body posted by a user to find a duplicate user question title from the English StackExchange forum
2728		
2729		
2730	CQADupstackGamingRetrieval Q	Represent the title of a user question to find a duplicate user question title with body from the Gaming StackExchange forum
2731		
2732	CQADupstackGamingRetrieval D	Represent the question title with body posted by a user to find a duplicate user question title from the Gaming StackExchange forum
2733		
2734		
2735		
2736	CQADupstackGisRetrieval Q	Represent the title of a user question to find a duplicate user question title with body from the Gis StackExchange forum
2737		
2738	CQADupstackGisRetrieval D	Represent the question title with body posted by a user to find a duplicate user question title from the Gis StackExchange forum
2739		
2740		
2741	CQADupstackMathematicaRetrieval Q	Represent the title of a user question to find a duplicate user question title with body from the Mathematica StackExchange forum
2742		
2743	CQADupstackMathematicaRetrieval D	Represent the question title with body posted by a user to find a duplicate user question title from the Mathematica StackExchange forum
2744		
2745		
2746		
2747	CQADupstackPhysicsRetrieval Q	Represent the title of a user question to find a duplicate user question title with body from the Physics StackExchange forum
2748		
2749	CQADupstackPhysicsRetrieval D	Represent the question title with body posted by a user to find a duplicate user question title from the Physics StackExchange forum
2750		
2751		
2752		
2753		

2754	CQADupstackProgrammersRetrieval Q	Represent the title of a user question to find a duplicate user question title with body from the Programmers StackExchange forum
2755		
2756		
2757	CQADupstackProgrammersRetrieval D	Represent the question title with body posted by a user to find a duplicate user question title from the Programmers StackExchange forum
2758		
2759		
2760		
2761	CQADupstackStatsRetrieval Q	Represent the title of a user question to find a duplicate user question title with body from the Stats StackExchange forum
2762		
2763	CQADupstackStatsRetrieval D	Represent the question title with body posted by a user to find a duplicate user question title from the Stats StackExchange forum
2764		
2765	CQADupstackTexRetrieval Q	Represent the title of a user question to find a duplicate user question title with body from the Tex StackExchange forum
2766		
2767	CQADupstackTexRetrieval D	Represent the question title with body posted by a user to find a duplicate user question title from the Tex StackExchange forum
2768		
2769		
2770	CQADupstackUnixRetrieval Q	Represent the title of a user question to find a duplicate user question title with body from the Unix StackExchange forum
2771		
2772	CQADupstackUnixRetrieval D	Represent the question title with body posted by a user to find a duplicate user question title from the Unix StackExchange forum
2773		
2774	CQADupstackWebmastersRetrieval Q	Represent the title of a user question to find a duplicate user question title with body from the Webmasters StackExchange forum
2775		
2776		
2777	CQADupstackWebmastersRetrieval D	Represent the question title with body posted by a user to find a duplicate user question title from the Webmasters StackExchange forum
2778		
2779		
2780		
2781	CQADupstackWordpressRetrieval Q	Represent the title of a user question to find a duplicate user question title with body from the Wordpress StackExchange forum
2782		
2783	CQADupstackWordpressRetrieval D	Represent the question title with body posted by a user to find a duplicate user question title from the Wordpress StackExchange forum
2784		
2785		
2786	DBPedia Q	Represent the entity to find a title with abstract about this entity from the DBPedia corpus
2787		
2788	DBPedia D	Represent the title with abstract of a DBPedia corpus entry to find the entity of a few words it is about
2789		
2790		
2791	FEVER Q	Represent the claim to find a Wikipedia abstract to support it
2792	FEVER D	Represent the Wikipedia abstract to find a claim that it supports
2793	FiQA2018 Q	Represent the StackExchange user query to find a StackExchange post from the Investment topic that answers it
2794		
2795	FiQA2018 D	Represent the StackExchange post from the Investment topic to find a StackExchange user query that it answers
2796		
2797		
2798	HotpotQA Q	Represent the multi-hop question to find a Wikipedia passage that answers it
2799		
2800	HotpotQA D	Represent the Wikipedia passage to find a multi-hop question that it answers
2801		
2802	MSMARCO Q	Represent the Bing user search query to find a passage that adequately addresses it
2803		
2804	MSMARCO D	Represent the passage for finding a Bing user search query about it
2805		
2806		
2807	NFCorpus Q	Represent the query from NutritionFacts to find a title with text of a medical document from PubMed about it

2808	NFCorpus D	Represent this text of a medical document from PubMed to find a query someone may enter at NutritionFacts that it answers
2809		
2810	NQ Q	Represent the Google search query to find an answer span from a Wikipedia article that addresses it
2811		
2812	NQ D	Represent the Wikipedia article span to find a Google search query that would be addressed by it
2813		
2814	SCIDOCS Q	Represent the scientific paper title to find the title with abstract of a scientific paper on PubMed that it has likely cited
2815		
2816	SCIDOCS D	Represent the title with abstract of this scientific paper to find the title of another scientific paper on PubMed that likely cites this article
2817		
2818	SciFact Q	Represent the scientific claim to find a scientific paper abstract from PubMed to support it
2819		
2820	SciFact D	Represent the scientific paper abstract from PubMed to find a scientific claim that it supports
2821		
2822	TRECCOVID Q	Represent the search query to find a scientific article about COVID-19 that adequately addresses the query
2823		
2824	TRECCOVID D	Represent the scientific article about COVID-19 to find a user query that it adequately addresses
2825		
2826	Touche2020 Q	Represent the question to find a title with passage of an argument from args.me that takes a stance about it
2827		
2828	Touche2020 D	Represent the title with passage of an argument from args.me to find a question that it takes a stance about
2829		
2830	QuoraRetrieval Q	Represent the Quora question to find another short duplicate question on Quora
2831		
2832	QuoraRetrieval D	Represent the Quora question to find another short duplicate question on Quora
2833		
2834	AskUbuntuDupQuestions Q	Represent the query to find a duplicate query on the AskUbuntu community forum
2835		
2836	AskUbuntuDupQuestions D	Represent the query to find a duplicate query on the AskUbuntu community forum
2837		
2838	MindSmallReranking Q	Represent the news headline to find another news headline that the same reader would enjoy
2839		
2840	MindSmallReranking D	Represent the news headline to find another news headline that the same reader would enjoy
2841		
2842	SciDocsRR Q	Represent the title to find a similar scientific paper title
2843		
2844	SciDocsRR D	Represent the title to find a similar scientific paper title
2845		
2846	StackOverflowDupQuestions Q	Represent the query to find a duplicate query on the StackOverflow Java/JavaScript/Python community forums
2847		
2848	StackOverflowDupQuestions D	Represent the query to find a duplicate query on the StackOverflow Java/JavaScript/Python community forums
2849		
2850	BIOSSES	Represent the text to find another biological statement with the same meaning
2851		
2852	SICK-R	Represent the sentence to find another sentence with the same meaning
2853		
2854	STS12	Represent the sentence to find another sentence with the same meaning
2855		
2856		
2857		
2858		
2859		
2860		
2861		

2862	STS13	Represent the sentence to find another sentence with the same meaning
2863		
2864	STS14	Represent the sentence to find another sentence with the same meaning
2865		
2866	STS15	Represent the sentence to find another sentence with the same meaning
2867		
2868	STS16	Represent the sentence to find another sentence with the same meaning
2869		
2870	STS17	Represent the sentence to find another sentence with the same meaning
2871		
2872	STS22	Represent the sentence to find another sentence with the same meaning
2873		
2874	STSBenchmark	Represent the sentence to find another sentence with the same meaning
2875		
2876	SummEval Q	Represent the human-written summary to find a high-quality machine-written summary of the same news article
2877		
2878	SummEval D	Represent the machine-written summary to find a human-written summary with similar quality of the same news article
2879		
2880		
2881		
2882		
2883		
2884		
2885		

## T.2 EMBEDDING FEW-SHOT PROMPTS

Table 33: **1-shot example for the model trained on E5S.** The example is appended to the respective instruction in Table 31 separated by two newlines.

Task Name	Instruction
Banking77Classification	For example given "I am still waiting on my card?", it would match with "card.arrival"
EmotionClassification	For example given "ive been feeling a little burdened lately wasnt sure why that was", it would match with "sadness"
ImdbClassification	For example given "If only to avoid making this type of film in the future. This film is interesting as an experiment but tells no cogent story.;br /;One might feel virtuous for sitting thru it because it touches on so many IMPORTANT issues but it does so without any discernable motive. The viewer comes away with no new perspectives (unless one comes up with one while one's mind wanders, as it will invariably do during this pointless film).;br /;One might better spend one's time staring out a window at a tree growing.;br /;", it would match with "negative"
BiorxivClusteringP2P	For example given "Association of CDH11 with ASD revealed by matched-gene co-expression analysis and mouse behavioral studies", it would match with "neuroscience"
TwitterSemEval2015	For example given "The Ending to 8 Mile is my fav part of the whole movie", it would match with "Those last 3 battles in 8 Mile are THE shit"

2916	TwitterURLCorpus	For example given "Liberals , dont let Donald Trump tarnish L.L. Beans sterling brand reputation ", it would match with "Liberals, Don&rsquo;t Let Donald Trump Tarnish L.L. Bean&rsquo;s Sterling Brand Reputation"
2917		
2918		
2919		
2920		
2921	SprintDuplicateQuestions	For example given "Why is it impossible for me to find a easy way to send a picture with text on my Kyocera DuraCore ?", it would match with "Send or receive a picture with text - Kyocera DuraCore"
2922		
2923		
2924		
2925	AskUbuntuDupQuestions	For example given "what is a short cut i can use to switch applications ?", you should retrieve "keyboard short cut for switching between two or more instances of the same application ?"
2926		
2927		
2928	ArguAna	For example given "People will die if we don't do animal testing Every year, 23 new drugs are introduced in the UK alone.[13] Almost all will be tested on animals. A new drug will be used for a long time. Think of all the people saved by the use of penicillin. If drugs cost more to test, that means drug companies will develop less. This means more people suffering and dying", you should retrieve "animals science science general ban animal testing junior Many of these drugs are "me too" drugs – ones with a slight change that doesn't make much difference to an existing drug. [14] So often the benefits from animal testing are marginal, and even if there was a slight increase in human suffering, it would be worth it based on the animal suffering saved."
2929		
2930		
2931		
2932		
2933		
2934		
2935		
2936		
2937		
2938		
2939	SCIDOCS	For example given "A Direct Search Method to solve Economic Dispatch Problem with Valve-Point Effect", you should retrieve "A Hybrid EP and SQP for Dynamic Economic Dispatch with Nonsmooth Fuel Cost Function Dynamic economic dispatch (DED) is one of the main functions of power generation operation and control. It determines the optimal settings of generator units with predicted load demand over a certain period of time. The objective is to operate an electric power system most economically while the system is operating within its security limits. This paper proposes a new hybrid methodology for solving DED. The proposed method is developed in such a way that a simple evolutionary programming (EP) is applied as a based level search, which can give a good direction to the optimal global region, and a local search sequential quadratic programming (SQP) is used as a fine tuning to determine the optimal solution at the final. Ten units test system with nonsmooth fuel cost function is used to illustrate the effectiveness of the proposed method compared with those obtained from EP and SQP alone."
2940		
2941		
2942		
2943		
2944		
2945		
2946		
2947		
2948		
2949	STS12	For example given "Counties with population declines will be Vermillion, Posey and Madison.", it would match with "Vermillion, Posey and Madison County populations will decline."
2950		
2951		
2952		
2953		
2954		
2955		
2956		
2957		
2958		
2959		
2960		
2961		
2962		
2963		
2964		
2965		
2966		
2967		
2968		
2969		

2970	SummEval	The provided query could be "Mexican restaurant has decided to tap into \$70 billion food delivery market. Fast-casual chain will work with the Postmates app to allow mobile orders. App works in similar way to Uber, using hired drivers to deliver the food. But the chain will add a 9% service charge - on top of Postmates\$5 rate." and the positive "chipotle has decided to tap into the \$ 70 billion food delivery market by teaming up with an app to bring burritos straight to customers doors . the fast-casual chain will work with the postmates app to begin offering delivery for online and mobile orders in 67 cities . the restaurant plans to add a nine per cent service charge - with the delivery fees for postmates beginning at \$ 5 and up depending on distance and demand ."
2971		
2972		
2973		
2974		
2975		
2976		
2977		
2978		
2979		
2980		
2981		
2982		

Table 34: **1-shot example for the model trained on MEDI2.** The example is appended to the respective instruction in Table 32 separated by two newlines.

Task Name	Instruction
Banking77Classification	The provided query could be "I am still waiting on my card?" and the positive "What can I do if my card still hasn't arrived after 2 weeks?"
EmotionClassification	The provided query could be "ive been feeling a little burdened lately wasnt sure why that was" and the positive "i feel like i have to make the suffering i m seeing mean something"
ImdbClassification	The provided query could be "If only to avoid making this type of film in the future. This film is interesting as an experiment but tells no cogent story.;br /;One might feel virtuous for sitting thru it because it touches on so many IMPORTANT issues but it does so without any discernable motive. The viewer comes away with no new perspectives (unless one comes up with one while one's mind wanders, as it will invariably do during this pointless film).;br /;One might better spend one's time staring out a window at a tree growing.;br /;" and the positive "The silent one-panel cartoon Henry comes to Fleischer Studios, billed as "The world's funniest human" in this dull little cartoon. Betty, long past her prime, thanks to the Production Code, is running a pet shop and leaves Henry in charge for far too long – five minutes. A bore."
SprintDuplicateQuestions	The provided query could be "Why is it impossible for me to find a easy way to send a picture with text on my Kyocera DuraCore ?" and the positive "Send or receive a picture with text - Kyocera DuraCore"
TwitterSemEval2015	For example given "The Ending to 8 Mile is my fav part of the whole movie", it would match with "Those last 3 battles in 8 Mile are THE shit"
TwitterURLCorpus	For example given "Liberals , dont let Donald Trump tarnish L.L. Beans sterling brand reputation ", it would match with "Liberals, Don&rsquo;t Let Donald Trump Tarnish L.L. Bean&rsquo;s Sterling Brand Reputation"



3024	AskUbuntuDupQuestions	The provided query could be "what is a short cut i can use to switch applications ?" and the positive "keyboard short cut for switching between two or more instances of the same application ?"
3025		
3026		
3027		
3028		
3029	ArguAna	The provided query could be "People will die if we don't do animal testing Every year, 23 new drugs are introduced in the UK alone.[13] Almost all will be tested on animals. A new drug will be used for a long time. Think of all the people saved by the use of penicillin. If drugs cost more to test, that means drug companies will develop less. This means more people suffering and dying" and the positive "animals science science general ban animal testing junior Many of these drugs are "me too" drugs – ones with a slight change that doesn't make much difference to an existing drug. [14] So often the benefits from animal testing are marginal, and even if there was a slight increase in human suffering, it would be worth it based on the animal suffering saved."
3030		
3031		
3032		
3033		
3034		
3035		
3036		
3037		
3038		
3039		
3040	SCIDOCS	The provided query could be "A Direct Search Method to solve Economic Dispatch Problem with Valve-Point Effect" and the positive "A Hybrid EP and SQP for Dynamic Economic Dispatch with Nonsmooth Fuel Cost Function Dynamic economic dispatch (DED) is one of the main functions of power generation operation and control. It determines the optimal settings of generator units with predicted load demand over a certain period of time. The objective is to operate an electric power system most economically while the system is operating within its security limits. This paper proposes a new hybrid methodology for solving DED. The proposed method is developed in such a way that a simple evolutionary programming (EP) is applied as a based level search, which can give a good direction to the optimal global region, and a local search sequential quadratic programming (SQP) is used as a fine tuning to determine the optimal solution at the final. Ten units test system with nonsmooth fuel cost function is used to illustrate the effectiveness of the proposed method compared with those obtained from EP and SQP alone."
3041		
3042		
3043		
3044		
3045		
3046		
3047		
3048		
3049		
3050		
3051		
3052		
3053		
3054		
3055		
3056		
3057		
3058	STS12	The provided query could be "Counties with population declines will be Vermillion, Posey and Madison." and the positive "Vermillion, Posey and Madison County populations will decline."
3059		
3060		
3061		
3062	SummEval	The provided query could be "Mexican restaurant has decided to tap into \$70 billion food delivery market. Fast-casual chain will work with the Postmates app to allow mobile orders. App works in similar way to Uber, using hired drivers to deliver the food. But the chain will add a 9% service charge - on top of Postmates\$5 rate." and the positive "chipotle has decided to tap into the \$ 70 billion food delivery market by teaming up with an app to bring burritos straight to customers doors . the fast-casual chain will work with the postmates app to begin offering delivery for online and mobile orders in 67 cities . the restaurant plans to add a nine per cent service charge - with the delivery fees for postmates beginning at \$ 5 and up depending on distance and demand ."
3063		
3064		
3065		
3066		
3067		
3068		
3069		
3070		
3071		
3072		
3073		
3074		
3075		
3076		
3077		

3078 T.3 GENERATIVE PROMPTS  
3079

3080 Figure 14 until Figure 19 contain the prompts with examples used for our generative tasks.  
3081

---

3082 **Input:**

---

3083  
3084 <s><|user|>

3085 The following are multiple choice questions (with answers) about abstract algebra.

3086 Find the degree for the given field extension  $Q(\sqrt{2}, \sqrt{3}, \sqrt{18})$  over  $Q$ .

3087 A. 0

3088 B. 4

3089 C. 2

3090 D. 6

3091 Answer:

3092 <|assistant|>

3093 The answer is:

---

3094  
3095 **Correct completion:**

---

3096  
3097 B

---

3098  
3099 Figure 14: MMLU prompt example.  
3100  
3101  
3102  
3103  
3104  
3105  
3106  
3107  
3108  
3109  
3110  
3111  
3112  
3113  
3114  
3115  
3116  
3117  
3118  
3119  
3120  
3121  
3122  
3123  
3124  
3125  
3126  
3127  
3128  
3129  
3130  
3131

3132  
3133  
3134  
3135  
3136  
3137  
3138  
3139  
3140  
3141  
3142  
3143  
3144  
3145  
3146  
3147  
3148  
3149  
3150  
3151  
3152  
3153  
3154  
3155  
3156  
3157  
3158  
3159  
3160  
3161  
3162  
3163  
3164  
3165  
3166  
3167  
3168  
3169  
3170  
3171  
3172  
3173  
3174  
3175  
3176  
3177  
3178  
3179  
3180  
3181  
3182  
3183  
3184  
3185

---

**Input:**

---

<s><|user|>

Answer the following questions.

Question: There are 15 trees in the grove. Grove workers will plant trees in the grove today. After they are done, there will be 21 trees. How many trees did the grove workers plant today?

Answer: There are 15 trees originally. Then there were 21 trees after some more were planted. So there must have been  $21 - 15 = 6$ . So the answer is 6.

Question: If there are 3 cars in the parking lot and 2 more cars arrive, how many cars are in the parking lot?

Answer: There are originally 3 cars. 2 more cars arrive.  $3 + 2 = 5$ . So the answer is 5.

Question: Leah had 32 chocolates and her sister had 42. If they ate 35, how many pieces do they have left in total?

Answer: Originally, Leah had 32 chocolates. Her sister had 42. So in total they had  $32 + 42 = 74$ . After eating 35, they had  $74 - 35 = 39$ . So the answer is 39.

Question: Jason had 20 lollipops. He gave Denny some lollipops. Now Jason has 12 lollipops. How many lollipops did Jason give to Denny?

Answer: Jason started with 20 lollipops. Then he had 12 after giving some to Denny. So he gave Denny  $20 - 12 = 8$ . So the answer is 8.

Question: Shawn has five toys. For Christmas, he got two toys each from his mom and dad. How many toys does he have now?

Answer: Shawn started with 5 toys. If he got 2 toys each from his mom and dad, then that is 4 more toys.  $5 + 4 = 9$ . So the answer is 9.

Question: There were nine computers in the server room. Five more computers were installed each day, from monday to thursday. How many computers are now in the server room?

Answer: There were originally 9 computers. For each of 4 days, 5 more computers were added. So  $5 * 4 = 20$  computers were added.  $9 + 20$  is 29. So the answer is 29.

Question: Michael had 58 golf balls. On tuesday, he lost 23 golf balls. On wednesday, he lost 2 more. How many golf balls did he have at the end of wednesday?

Answer: Michael started with 58 golf balls. After losing 23 on tuesday, he had  $58 - 23 = 35$ . After losing 2 more, he had  $35 - 2 = 33$  golf balls. So the answer is 33.

Question: Olivia has \$23. She bought five bagels for \$3 each. How much money does she have left?

Answer: Olivia had 23 dollars. 5 bagels for 3 dollars each will be  $5 * 3 = 15$  dollars. So she has  $23 - 15$  dollars left.  $23 - 15$  is 8. So the answer is 8.

Question: The girls are trying to raise money for a carnival. Kim raises \$320 more than Alexandra, who raises \$430, and Maryam raises \$400 more than Sarah, who raises \$300. How much money, in dollars, did they all raise in total?

<|assistant|>

Answer:

---

**Correct completion:**

---

Kim raises  $320+430=750$  dollars. Maryam raises  $400+300=700$  dollars. They raise  $750+430+400+700=2280$  dollars. So the answer is 2280.

---

Figure 15: **GSM8K prompt example.**

3186  
3187  
3188  
3189 **Input:**  
3190  
3191 <s><|user|>  
3192 Questions that involve enumerating objects and asking the model to count them.  
3193  
3194 Q: I have a blackberry, a clarinet, a nectarine, a plum, a strawberry, a banana, a flute, an  
3195 orange, and a violin. How many fruits do I have?  
3196 A: Let's think step by step.  
3197 We first identify the fruits on the list and include their quantity in parentheses:  
3198 - blackberry (1)  
3199 - nectarine (1)  
3200 - plum (1)  
3201 - strawberry (1)  
3202 - banana (1)  
3203 - orange (1)  
3204 Now, let's add the numbers in parentheses:  $1 + 1 + 1 + 1 + 1 + 1 = 6$ . So the answer is 6.  
3205  
3206 Q: I have an orange, a raspberry, two peaches, a blackberry, an apple, a grape, a nectarine, and three  
3207 plums. How many fruits do I have?  
3208 A: Let's think step by step.  
3209 We first identify the fruits on the list and include their quantity in parentheses:  
3210 - orange (1)  
3211 - raspberry (1)  
3212 - peaches (2)  
3213 - blackberry (1)  
3214 - apple (1)  
3215 - grape (1)  
3216 - nectarine (1)  
3217 - plums (3)  
3218 Now, let's add the numbers in parentheses:  $1 + 1 + 2 + 1 + 1 + 1 + 1 + 3 = 11$ . So the answer is 11.  
3219  
3220 Q: I have a lettuce head, a head of broccoli, an onion, a stalk of celery, two carrots, a garlic, and a  
3221 yam. How many vegetables do I have?  
3222 A: Let's think step by step.  
3223 We first identify the vegetables on the list and include their quantity in parentheses:  
3224 - lettuce (1)  
3225 - broccoli (1)  
3226 - onion (1)  
3227 - celery (1)  
3228 - carrots (2)  
3229 - garlic (1)  
3230 - yam (1)  
3231 Now, let's add the numbers in parentheses:  $1 + 1 + 1 + 1 + 2 + 1 + 1 = 8$ . So the answer is 8.  
3232  
3233 Q: I have a banana, four strawberries, an apple, two peaches, a plum, a blackberry, and two raspberries.  
3234 How many fruits do I have?  
3235 <|assistant|>  
3236  
3237 **Correct completion:**  
3238  
3239 12

Figure 16: **BBH prompt example.**

3240  
3241  
3242  
3243  
3244  
3245  
3246  
3247  
3248  
3249  
3250  
3251  
3252  
3253  
3254  
3255  
3256  
3257  
3258  
3259  
3260  
3261  
3262  
3263  
3264  
3265  
3266  
3267  
3268  
3269  
3270  
3271  
3272  
3273  
3274  
3275  
3276  
3277  
3278  
3279  
3280  
3281  
3282  
3283  
3284  
3285  
3286  
3287  
3288  
3289  
3290  
3291  
3292  
3293

---

**Input:**

---

<s><|user|>

Jawab pertanyaan berikut berdasarkan informasi di bagian yang diberikan.

Bagian: Mula-mula pada pelukis seorang pelukis pemandangan Wahdi Sumanta, Abdullah Suriosubroto (ayah Basuki Abdullah). Kemudian bertemu dan berkenalan dengan Affandi, Sudarso, dan Barli. Mereka lalu membentuk kelompok Lima serangkai. Di rumah tempat tinggal Affandi mereka mengadakan latihan melukis bersama dengan tekun dan mendalam. Dari Wahdi, ia banyak menggali pengetahuan tentang melukis. Kegiatannya bukan hanya melukis semata, tetapi pada waktu senggang ia menceburkan diri pada kelompok sandiwara Sunda sebagai pelukis dekor. Dari pengalaman itulah, ia mengasah kemampuannya.

Pertanyaan: dari manakah Hendra Gunawan belajar melukis?

Jawaban: kelompok Lima serangkai

Bagian: Empat Sehat Lima Sempurna adalah kampanye yang dilakukan pemerintah sejak tahun 1955 untuk membuat masyarakat memahami pola makan yang benar.[1]. Dalam konsep 4 sehat 5 sempurna, makanan dibagi atas empat sumber nutrisi penting, yaitu makanan pokok, lauk pauk, sayur-mayur, buah-buahan, dan disempurnakan dengan susu bila mampu, menjadi lima sempurna[2] Konsep ini menekankan pentingnya empat golongan makanan berupa sumber kalori untuk tenaga, protein untuk pembangunan, sayur dan buah sumber vitamin dan mineral untuk pemeliharaan.[1]

Pertanyaan: siapakah yang menciptakan Ide 4 sehat 5 sempurna pertama kali?

<|assistant|>

Jawaban:

---

**Correct completion:**

---

pemerintah

---

Figure 17: TyDi QA prompt example from Indonesian.

3294  
3295  
3296  
3297  
3298  
3299  
3300  
3301  
3302  
3303  
3304  
3305  
3306  
3307  
3308  
3309  
3310  
3311  
3312  
3313  
3314  
3315  
3316  
3317  
3318  
3319  
3320  
3321  
3322  
3323  
3324  
3325  
3326  
3327  
3328  
3329  
3330  
3331  
3332  
3333  
3334  
3335  
3336  
3337  
3338  
3339  
3340  
3341  
3342  
3343  
3344  
3345  
3346  
3347

---

**Input:**

---

```
<|user|>
Write a Python function 'has_close_elements(numbers: List[float], threshold: float) -> bool' to solve
the following problem: Check if in given list of numbers, are any two numbers closer to each other
than given threshold.

>>> has_close_elements([1.0, 2.0, 3.0], 0.5)
False
>>> has_close_elements([1.0, 2.8, 3.0, 4.0, 5.0, 2.0], 0.3)
True

<|assistant|>

from typing import List
def has_close_elements(numbers: List[float], threshold: float) -> bool:
    """ Check if in given list of numbers, are any two numbers closer to
        each other than
        given threshold.
    """
    >>> has_close_elements([1.0, 2.0, 3.0], 0.5)
    False
    >>> has_close_elements([1.0, 2.8, 3.0, 4.0, 5.0, 2.0], 0.3)
    True
    """
```

---

**Correct completion:**

---

```
for idx, elem in enumerate(numbers):
    for idx2, elem2 in enumerate(numbers):
        if idx != idx2:
            distance = abs(elem - elem2)
            if distance < threshold:
                return True

return False
```

---

Figure 18: **HumanEvalSynthesize** prompt example.

3348  
3349  
3350  
3351  
3352  
3353  
3354  
3355  
3356  
3357  
3358  
3359  
3360  
3361  
3362  
3363  
3364  
3365  
3366  
3367  
3368  
3369  
3370  
3371  
3372  
3373  
3374  
3375  
3376  
3377  
3378  
3379  
3380  
3381  
3382  
3383  
3384  
3385  
3386  
3387  
3388  
3389  
3390  
3391  
3392  
3393  
3394  
3395  
3396  
3397  
3398  
3399  
3400  
3401

---

**Input:**

---

<s><|user|>

If you were a Shakespearean character, how would you declare your love for someone in a soliloquy?

---

**Completion by GPT-4 Turbo:**

---

As an AI, I lack personal emotions and experiences, but I can craft a soliloquy in the style of a Shakespearean character declaring love. Here's an example:

(Imagine a dimly lit stage, a single spotlight shining down upon the forlorn lover, who steps forward to share the depths of their heart.)

O, speak, thou lustrous beacon of the night,  
With stars enkindled in thine eyes so bright.  
Thou art the daybreak to my endless gloom,  
The rose that in my barren heart doth bloom.

Forsooth, this love, a tempest full of grace,  
Hath swept upon my soul without a trace.  
My heart, once caged, now beats beneath thy gaze,  
And yearns to laud thy beauty all its days.

What art thou, Love? A tyrant or a muse?  
That rends my sense and all my thoughts confuse?  
Thou art the riddle none can e'er resolve,  
Yet in thy puzzle, my life doth revolve.

With every breath, I thee do magnify,  
And to the world, my passion testify.  
Thine essence hath become my sole refrain,  
A melody that in my mind doth reign.

Shall I compare thee to a summer's day?  
Thou art more lovely and more temperate.  
Rough winds do shake the darling buds of May,  
And summer's lease hath all too short a date.

But thy eternal summer shall not fade,  
Nor lose possession of that fair thou ow'st.  
Nor shall Death brag thou wander'st in his shade,  
When in eternal lines to time thou grow'st.

So long as men can breathe, or eyes can see,  
So long lives this, and this gives life to thee.

(With a heavy sigh, the lover steps back, their heart laid bare, the words hanging in the air like the fragrance of an unseen bloom, their love declared, though the object of their affection remains unaware, the audience enthralled by the vulnerability of true, unspoken love.)

---

Figure 19: AlpacaEval prompt example.

3402 T.4 RAG PROMPTS

3403

3404

---

**Generative input:**

3405

3406

<s><|user|>

3407

what was the name of darth vader star destroyer

3408

<|assistant|>

3409

The answer is

3410

---

**Generation:**

3411

3412

not clear. Darth Vader had a star destroyer, but it was

3413

3414

---

Figure 20: No RAG example and generation by GRITLM 7B.

3415

3416

3417

---

**Generative input:**

3418

3419

<s><|user|>

3420

what was the name of darth vader star destroyer

3421

3422

Star Destroyer The iconic Star Destroyer first appears in the opening scene of Star Wars, as Darth Vader’s flagship, the Devastator, chases the Tantive IV above Tatooine. This appearance shows the Imperial ship’s massive size in comparison to the Tantive IV.

3423

3424

3425

Optionally using the prior context answer the query prior to it

3426

<|assistant|>

3427

The answer is

3428

3429

---

**Generation:**

3430

3431

the Devastator.</s>

3432

3433

---

Figure 21: RAG query then document example and generation by GRITLM 7B.

3434

3435

3436

3437

3438

3439

3440

3441

3442

3443

3444

3445

3446

3447

3448

3449

3450

3451

3452

3453

3454

3455



3456  
3457  
3458  
3459  
3460  
3461  
3462  
3463  
3464  
3465  
3466  
3467  
3468  
3469  
3470  
3471  
3472  
3473  
3474  
3475  
3476  
3477  
3478  
3479  
3480  
3481  
3482  
3483  
3484  
3485  
3486  
3487  
3488  
3489  
3490  
3491  
3492  
3493  
3494  
3495  
3496  
3497  
3498  
3499  
3500  
3501  
3502  
3503  
3504  
3505  
3506  
3507  
3508  
3509

---

**Generative input:**

---

<s><|user|>  
Star Destroyer The iconic Star Destroyer first appears in the opening scene of Star Wars, as Darth Vader's flagship, the Devastator, chases the Tantive IV above Tatooine. This appearance shows the Imperial ship's massive size in comparison to the Tantive IV.

what was the name of darth vader star destroyer

Answer the prior query while optionally using the context prior to it  
<|assistant|>  
The answer is

---

**Generation:**

---

the Devastator.</s>

---

Figure 22: RAG document then query example and generation by GRITLM 7B.

---

**Embedding input (query) passed via key-value states:**

---

<s><|embed|>  
what was the name of darth vader star destroyer

---

**Generative input:**

---

\n<|user|>  
Star Destroyer The iconic Star Destroyer first appears in the opening scene of Star Wars, as Darth Vader's flagship, the Devastator, chases the Tantive IV above Tatooine. This appearance shows the Imperial ship's massive size in comparison to the Tantive IV.

Optionally using the prior context answer the query prior to it  
<|assistant|>  
The answer is

---

**Generation:**

---

Star Destroyer.</s>

---

Figure 23: GRIT Query Caching example and generation by GRITLM 7B.

3510  
3511  
3512 **Embedding input (doc) passed via key-value states and cached in the index:**  
3513  
3514 `<s><|embed|>`  
3515 Star Destroyer The iconic Star Destroyer first appears in the opening scene of Star Wars, as Darth  
3516 Vader’s flagship, the Devastator, chases the Tantive IV above Tatooine. This appearance shows the  
3517 Imperial ship’s massive size in comparison to the Tantive IV.  
3518  
3519 **Generative input:**  
3520  
3521 `\n<|user|>`  
3522 what was the name of darth vader star destroyer  
3523  
3524 Answer the prior query while optionally using the context prior to it  
3525 `<|assistant|>`  
3526 The answer is  
3527  
3528 **Generation:**  
3529  
3530 Devastator. The iconic Star Destroyer first appears in the opening

3531 Figure 24: **GRIT Doc Caching example and generation by GRITLM 7B.**  
3532  
3533  
3534  
3535  
3536

3537  
3538 **Embedding input (doc) passed via key-value states and cached in the index:**  
3539  
3540 `<s><|embed|>`  
3541 Star Destroyer The iconic Star Destroyer first appears in the opening scene of Star Wars, as Darth  
3542 Vader’s flagship, the Devastator, chases the Tantive IV above Tatooine. This appearance shows the  
3543 Imperial ship’s massive size in comparison to the Tantive IV.  
3544  
3545 **Embedding input (query) passed via key-value states:**  
3546  
3547 `<s><|embed|>`  
3548 what was the name of darth vader star destroyer  
3549  
3550 **Generative input:**  
3551  
3552 `\n<|user|>`  
3553 Answer the prior query while optionally using the context prior to it  
3554 `<|assistant|>`  
3555 The answer is  
3556  
3557 **Generation:**  
3558  
3559 the Star Destroyer. The Star Destroyer is a massive spacecraft

3559 Figure 25: **GRIT Doc-Query Caching example and generation by GRITLM 7B.** Unlike for Doc  
3560 Caching, we prepend the bos token (“< s >”) to both query and document, which improved the  
3561 match score from 14.13 to 18.39.  
3562  
3563

3564	<b>Embedding input (query) passed via key-value states:</b>
3565	
3566	<s>< embed >
3567	what was the name of darth vader star destroyer
3568	
3569	<b>Embedding input (doc) passed via key-value states and cached in the index:</b>
3570	
3571	< embed >
3572	Star Destroyer The iconic Star Destroyer first appears in the opening scene of Star Wars, as Darth
3573	Vader’s flagship, the Devastator, chases the Tantive IV above Tatooine. This appearance shows the
3574	Imperial ship’s massive size in comparison to the Tantive IV.
3575	<b>Generative Input:</b>
3576	
3577	\n< user >
3578	Optionally using the prior context answer the query prior to it
3579	< assistant >
3580	The answer is
3581	
3582	<b>Generation:</b>
3583	
3584	the Star Destroyer.

Figure 26: GRIT Query-Doc Caching example and generation by GRITLM 7B.

## U HARDWARE

For the training of GRITLM 7B, we used 8 nodes with 8 NVIDIA A100 80GB GPUs each for 48 hours corresponding to 3,072 GPU hours. Meanwhile for GRITLM 8x7B, we used 32 nodes with 8 NVIDIA H100 80GB GPUs each for 80 hours corresponding to 20,480 GPU hours. As we train both models for 1253 steps, this corresponds to several minutes per step. This slow training time is mainly due to (a) a large batch size per step, (b) large models and our associated strategies to make them fit into memory at the cost of speed (Appendix K, Appendix L), and (c) a cluster with slow inter-node communication. The Gen.-only and Emb.-only models in Table 1 used 72 and 1760 H100 80GB GPU hours, respectively. Adding up all ablations and evaluations, we likely used somewhere around 100,000 GPU hours.

## V ARTIFACTS

Table 35: Produced artifacts that will be released upon deanonimization.

Artifact	Public Link
<i>Table 6</i>	
7B KTO	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
8x7B KTO	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
<i>Table 10</i>	
CCCC WM	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
CCCC LT	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
BBCC M	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>

3618  
 3619  
 3620  
 3621  
 3622  
 3623  
 3624  
 3625  
 3626  
 3627  
 3628  
 3629  
 3630  
 3631  
 3632  
 3633  
 3634  
 3635  
 3636  
 3637  
 3638  
 3639  
 3640  
 3641  
 3642  
 3643  
 3644  
 3645  
 3646  
 3647  
 3648  
 3649  
 3650  
 3651  
 3652  
 3653  
 3654  
 3655  
 3656  
 3657  
 3658  
 3659  
 3660  
 3661  
 3662  
 3663  
 3664  
 3665  
 3666  
 3667  
 3668  
 3669  
 3670  
 3671

*Table 11*


---

CC WM	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
CB M	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
BB M	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>

---

*Table 12*


---

CC	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
BC	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
BC IL	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>

---

*Table 13*


---

Mistral 7B	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
Llama 2 7B	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
GPT-J 6B	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>

---

*Table 14*


---

MEDI	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
MEDI2 NNI	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
MEDI2	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
MEDI2 + W	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>

---

*Table 15*


---

MEDI	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
MEDI2	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
BBCC MEDI	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
BBCC MEDI2	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
BBCC MEDI2BGE	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
BBCC E5	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>

---

*Table 16*


---

Tülu 2 1 EP	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
Tülu 2 2 EP	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
OASST 1 EP	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
OASST 2 EP	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
UltraChat	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>

---

*Table 17*


---

No head	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
-> 1024	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>

---

*Table 18*


---

256	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
-----	---

---

3672	4096	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3673		
3674	<i>Table 19</i>	
3675		
3676	BF16	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3677	FP32	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3678		
3679	<i>Table 20</i>	
3680		
3681	Any dataset	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3682	Same dataset	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3683		
3684	<i>Table 21</i>	
3685		
3686	Tülu 2	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3687	Zephyr $\beta$	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3688		
3689	<i>Table 22</i>	
3690		
3691	MEDI 2048	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3692	MEDI 4096	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3693	BBCC MEDI2 512	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3694	BBCC MEDI2 2048	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3695		
3696	<i>Table 23</i>	
3697		
3698	E5 Token 4.2	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3699	E5 Token 6.0	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3700	E5 Mix 32 -> 8	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3701	E5 Mix 32 -> 8	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3702	MEDI2 Mix 4 -> 64	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3703	MEDI2 Mix 32 -> 8	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3704		
3705	<i>Other</i>	
3706		
3707	Code	<a href="https://github.com/ANONYMIZED">https://github.com/ANONYMIZED</a>
3708	Logs	<a href="https://wandb.ai/ANONYMIZED">https://wandb.ai/ANONYMIZED</a>
3709	Tülu 2	<a href="https://hf.co/datasets/ANONYMIZED">https://hf.co/datasets/ANONYMIZED</a>
3710	MEDI	<a href="https://hf.co/datasets/ANONYMIZED">https://hf.co/datasets/ANONYMIZED</a>
3711	MEDI	<a href="https://hf.co/datasets/ANONYMIZED">https://hf.co/datasets/ANONYMIZED</a>
3712	MEDI2	<a href="https://hf.co/datasets/ANONYMIZED">https://hf.co/datasets/ANONYMIZED</a>
3713	MEDI2	<a href="https://hf.co/datasets/ANONYMIZED">https://hf.co/datasets/ANONYMIZED</a>
3714	MEDI2BGE	<a href="https://hf.co/datasets/ANONYMIZED">https://hf.co/datasets/ANONYMIZED</a>
3715	GRITLM 7B NQ In-	<a href="https://hf.co/datasets/ANONYMIZED">https://hf.co/datasets/ANONYMIZED</a>
3716	dex (§5)	
3717		
3718	<i>Main artifacts</i>	
3719		
3720	GRITLM 7B	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3721	GRITLM 8x7B	<a href="https://hf.co/ANONYMIZED">https://hf.co/ANONYMIZED</a>
3722		
3723		
3724		
3725		

Table 36: **Used artifacts released by others.**

Model / Dataset	Public Link
GPT-4 (OpenAI et al., 2023)	<a href="https://openai.com/gpt-4">https://openai.com/gpt-4</a>
OpenAI v3 (OpenAI et al., 2023)	<a href="https://openai.com/blog/new-embedding-models-and-api-updates">https://openai.com/blog/new-embedding-models-and-api-updates</a>
Gemini (Team et al., 2023)	<a href="https://deepmind.google/technologies/gemini/">https://deepmind.google/technologies/gemini/</a>
Llama 2 (Touvron et al., 2023)	<a href="https://hf.co/meta-llama">https://hf.co/meta-llama</a>
Mistral 7B (Jiang et al., 2023)	<a href="https://hf.co/mistralai/Mistral-7B-v0.1">https://hf.co/mistralai/Mistral-7B-v0.1</a>
Mistral 7B Instruct (Jiang et al., 2023)	<a href="https://hf.co/mistralai/Mistral-7B-Instruct-v0.1">https://hf.co/mistralai/Mistral-7B-Instruct-v0.1</a>
Mixtral 8x7B (Jiang et al., 2024)	<a href="https://hf.co/mistralai/Mixtral-8x7B-v0.1">https://hf.co/mistralai/Mixtral-8x7B-v0.1</a>
Mixtral 8x7B Instruct (Jiang et al., 2024)	<a href="https://hf.co/mistralai/Mixtral-8x7B-Instruct-v0.1">https://hf.co/mistralai/Mixtral-8x7B-Instruct-v0.1</a>
Tülu 2 (Iverson et al., 2023)	<a href="https://hf.co/collections/allenai/tulu-v2-suite-6551b56e743e6349aab45101">https://hf.co/collections/allenai/tulu-v2-suite-6551b56e743e6349aab45101</a>
GPT-J 6B (Wang & Komatsuzaki, 2021)	<a href="https://hf.co/EleutherAI/gpt-j-6b">https://hf.co/EleutherAI/gpt-j-6b</a>
SGPT BE 5.8B (Muennighoff, 2022)	<a href="https://hf.co/Muennighoff/SGPT-5.8B-weightedmean-msmarco-specb-bitfit">https://hf.co/Muennighoff/SGPT-5.8B-weightedmean-msmarco-specb-bitfit</a>
Instructor-XL 1.5B (Su et al., 2023)	<a href="https://hf.co/hkunlp/instructor-xl">https://hf.co/hkunlp/instructor-xl</a>
BGE Large 0.34B (Xiao et al., 2023)	<a href="https://hf.co/BAAI/bge-large-en-v1.5">https://hf.co/BAAI/bge-large-en-v1.5</a>
Zephyr 7B $\beta$ (Tunstall et al., 2023)	<a href="https://hf.co/HuggingFaceH4/zephyr-7b-beta">https://hf.co/HuggingFaceH4/zephyr-7b-beta</a>
E5 Mistral 7B (Wang et al., 2024)	<a href="https://hf.co/intfloat/e5-mistral-7b-instruct">https://hf.co/intfloat/e5-mistral-7b-instruct</a>
UltraChat (Ding et al., 2023; Tunstall et al., 2023)	<a href="https://hf.co/datasets/HuggingFaceH4/ultrachat_200k">https://hf.co/datasets/HuggingFaceH4/ultrachat_200k</a>
OASST (Köpf et al., 2023; Muennighoff et al., 2023a)	<a href="https://hf.co/datasets/bigcode/oasst-octopack">https://hf.co/datasets/bigcode/oasst-octopack</a>